

Spatially-Dependent Multiple Testing Under Model Misspecification, with Application to Detection of Anthropogenic Influence on Extreme Climate Events

Mark D. Risser *

Climate and Ecosystem Sciences Division, Lawrence Berkeley National Laboratory
and

Christopher J. Paciorek

Department of Statistics, University of California, Berkeley
and

Dáithí A. Stone

Computational Research Division, Lawrence Berkeley National Laboratory

March 30, 2017

Abstract

The Weather Risk Attribution Forecast (WRAF) is a forecasting tool that uses output from global climate models to make simultaneous attribution statements about whether and how greenhouse gas emissions have contributed to extreme weather across the globe. However, in conducting a large number of simultaneous hypothesis tests, the WRAF is prone to identifying false “discoveries.” A common technique for addressing this multiple testing problem is to adjust the procedure in a way that controls the proportion of true null hypotheses that are incorrectly rejected, or the false discovery rate (FDR). Unfortunately, generic FDR procedures suffer from low power when the hypotheses are dependent, and techniques designed to account for dependence are sensitive to misspecification of the underlying statistical model. In this paper, we develop a Bayesian decision theoretic approach for dependent multiple testing that flexibly controls false discovery and is robust to model misspecification. We illustrate the robustness of our procedure to model error with a simulation study, using a framework that accounts for generic spatial dependence and allows the practitioner to flexibly specify the decision criteria. Finally, we outline the best procedure of those considered for use in the WRAF workflow and apply the procedure to several seasonal forecasts.

Keywords: False Discovery Rate, Decision Theory, Event Attribution, Climate Models, Empirical Orthogonal Functions

*The authors gratefully acknowledge *the Office of Science, Office of Biological and Environmental Research of the U.S. Department of Energy.*

1 Introduction

Event attribution (EA) is a field of study that seeks to understand and describe the influence of greenhouse gas emissions and other human activities on extreme weather (Stott et al., 2013; National Academies of Sciences and Medicine, 2016). The increasing interest in this field arises from the realization that a major fraction of past, current, and future climate impacts and climate change-related impacts result from the occurrence of extreme weather (Arent et al., 2014; Smith et al., 2014). Risk-based EA studies quantify the effect of greenhouse gas (GHG) emissions and other anthropogenic factors on weather by comparing two climate scenarios: a factual real-world scenario (the “world as it is”) and a counterfactual, non-anthropogenic world (the “world as it might have been”). Then, using a probabilistic framework (Allen, 2003; Stone and Allen, 2005; Hansen et al., 2014), a risk-based EA study compares the probabilities of pre-defined unusual weather in the two scenarios and estimates how much more or less likely extreme events are in the anthropogenically-influenced world than they would have been otherwise (note: here and throughout we mean “risk” in the sense of epidemiological or relative risk, not statistical risk). Typically, the probabilities for each of these scenarios are estimated from simulations of climate models.

Risk-based EA studies can either be targeted or systematic in their approach. Targeted studies examine one event (or a small number of events), studying in detail the meteorological mechanisms involved in the event and how the anthropogenic influence is transmitted through them, and are generally reactive in the sense that they are only conducted for an event that has actually occurred (e.g. Stott et al., 2004; Pall et al., 2011). Systematic studies cover a much larger number of events using an identical method for all events, but the rigidity of a single experimental design means that some events are not amenable to investigation (Angélil et al., 2017). An advantage of the systematic approach is that it does not necessarily depend on the occurrence of events, with it being possible to instead perform the analyses on a pre-defined list of events. This is the approach taken by the Weather Risk Attribution Forecast (WRAF, <http://www.csag.uct.ac.za/~daithi/forecast>). In order to have EA information available in “real-time”, the WRAF performs analyses one month in advance using a pre-defined list of 232 extreme weather events, comprising an unusually hot, cold, wet, and/or dry month over each of 58 regions (Angélil et al., 2014). In the upcoming new version, the number of regions will be increased by a factor of about four (see, e.g., Figure 4; D. Stone, “A hierarchical collection of political/economic regions

for analysis of climate extremes”, in preparation). Data for both the factual and counterfactual scenarios come from climate model simulations.

Formally, the forecast involves estimating the probability of a pre-defined extreme event for both climate scenarios in each of the regions. For region $i = 1, \dots, M$ (in the upcoming version of the WRAF, $M = 237$), the forecast uses the ratio of scenario-specific probabilities p_{Fi} (for the factual scenario) and p_{Ci} (for the counterfactual scenario) or “risk ratio” $RR_i = p_{Fi}/p_{Ci}$ to formally test for changes in the probability of an extreme month. In other words, a collection of statistical tests are conducted that have null hypotheses of the form

$$H_i : RR_i \leq c, \quad i = 1, \dots, M, \quad (1)$$

where, for example, $c = 1$ if we are interested in determining whether anthropogenic influence has resulted in an increase in the event probability. Ultimately, we wish to separately test collections of hypotheses like (1) for extreme temperature (both hot and cold) and precipitation (both wet and dry).

Of course, when the number of tests M is large, a classical testing procedure is prone to identifying false “discoveries,” or incorrectly rejecting null hypotheses (commonly referred to as Type I errors). As such, the testing procedure is often adjusted, attempting to control the false discovery rate (FDR), which is the proportion of true null hypotheses that are incorrectly rejected. Since the data arise from physical climate models, it is anticipated that the hypotheses might be dependent: in other words, there is likely strong dependence within each spatial field of probabilities. This dependence might arise from the spatial proximity of the regions (i.e., strong dependence between p_{Fi} and p_{Fj} for adjacent regions i and j) but also from potentially un-specified long-range teleconnections (in which two probabilities p_{Fi} and p_{Fj} might be highly correlated even if regions i and j are far apart) that are common for atmospheric climate variables considered over the globe (see, e.g., [Cressie and Wikle, 2011](#)). Unfortunately, while classical FDR procedures ([Benjamini and Hochberg, 1995](#); see Section 2) are theoretically valid for positively correlated hypotheses ([Benjamini and Yekutieli, 2001](#)), they are also known to suffer from low power when the test statistics from each test are not independent (e.g., see [Sun and Cai, 2009](#)). And, while the literature contains a number of methods for applying FDR procedures under dependence, the methods are outlined for specific underlying probability models and are sensitive to improper specification of this model

([Sun et al., 2015](#)).

In this paper, we develop an approach to the multiple testing problem for spatially-dependent hypotheses in a systematic and decision-theoretic framework. Focusing on procedures that account for dependence among tests, we provide an overview of the diverse literature on false discovery control (including traditional methods and both Frequentist and Bayesian decision-theoretic approaches). The framework we use (originally introduced by [Müller et al., 2004](#)) allows the practitioner to flexibly specify the decision criteria for false discovery control, and we explore practical comparison of various FDR procedures and decision criteria. When an empirical estimate of the correlation among tests is available, we introduce a robust yet practical modeling framework for addressing spatial dependence among hypotheses. Furthermore, we specifically address sensitivity of the performance of the decision rule to statistical model misspecification and demonstrate the robustness of the modeling framework for FDR control. While the methodology is designed specifically for the hypothesis testing setting of the WRAF, our framework is useful for a broader set of problems involving multiple testing over a spatial domain, particularly in the case where an empirical correlation estimate is available, which is often the case for climate science scenarios. In this context, the methods outlined in this paper could be used for general hierarchical Bayesian models beyond just considering the probability of extremes or the risk ratio.

A reader familiar with the climate science literature will be aware of the concept of statistical field significance ([Livezey and Chen, 1983](#)), which is an alternative multiple testing approach that seeks to evaluate the collective significance of a set of statistics. While field significance techniques are well-established in climate science, we instead seek to control FDR following the arguments outlined in [Ventura et al. \(2004\)](#), the most important of which is that field significance provides no specific information about which individual tests are significant. Interestingly, the idea of FDR-control is growing in popularity among climate scientists, as evidenced by a recent paper by [Wilks \(2016\)](#). Finally, note that if one is only concerned with a real-time attribution statement for a single region in advance, then the multiple testing framework presented here is not required.

The paper proceeds as follows. In [Section 2](#), we introduce a decision theoretic framework for FDR control and present a systematic and flexible Bayesian approach to the problem. In [Section 3](#), we conduct a simulation study to assess the sensitivity of the Bayesian procedure to misspecification of the statistical model and identify a data-driven approach that robustly controls the FDR. In [Section 4](#), we apply the method to a real data set to be used for the WRAF; [Section 5](#) concludes

the paper.

2 Decision theoretic approaches for false discovery control

Formally, define scenario-specific data $\mathbf{Z} = \{(Z_{Fi}, Z_{Ci}) : i = 1, \dots, M\}$ (see Section 3) and a set of null hypotheses $\{H_i : i = 1, \dots, M\}$, as well as a corresponding collection of unknown parameters that represent the true state for each hypothesis:

$$\theta_i = \begin{cases} 0 & \text{if the true state of hypothesis } i \text{ is null} \\ 1 & \text{if the true state of hypothesis } i \text{ is non-null} \end{cases} \quad i = 1, \dots, M.$$

The testing problem involves generating a decision rule $\delta(\mathbf{Z}) \equiv \boldsymbol{\delta} = \{\delta_i : i = 1, \dots, M\}$, such that

$$\delta_i = \begin{cases} 0 & \text{if hypothesis } i \text{ is classified as null} \\ 1 & \text{if hypothesis } i \text{ is classified as non-null} \end{cases} \quad i = 1, \dots, M.$$

In addition to specifying the form of the decision rule (often based on a test statistic, P -value, etc.), an underlying probability model must be specified in order to estimate the decision rule. The false discovery proportion (FDP) is defined as $\text{FDP} = [\sum_{i=1}^M (1 - \theta_i) \delta_i] / [\sum_{i=1}^M \delta_i]$. Note that the FDP is simply a function of unknown parameters (θ_i) and random variables (δ_i), and is hence fundamentally neither Frequentist nor Bayesian.

The original FDR procedure given by [Benjamini and Hochberg \(1995\)](#) (henceforth BH) controls the (frequentist) FDR, defined as the expected FDP, i.e., $\text{FDR} \equiv \text{E}(\text{FDP})$, where the expectation is taken over repeated experiments. Their remarkably simple procedure ensures that $\text{FDR} \leq \alpha$; the proof in [Benjamini and Hochberg \(1995\)](#) is established for independent test statistics and any configuration of false null hypotheses. Alternatives to the original procedure include an adaptive FDR procedure ([Benjamini and Hochberg, 2000](#), [Genovese and Wasserman, 2002](#); henceforth AP); procedures that control either the positive FDR, $\text{pFDR} = \text{E}(\text{FDP} \mid \sum_{i=1}^M \delta_i > 0)$ or the marginal FDR, $\text{mFDR} = \text{E}(\sum_{i=1}^M (1 - \theta_i) \delta_i) / \text{E}(\sum_{i=1}^M \delta_i)$ ([Storey, 2003](#)); and Bayesian approaches to the problem using local FDR ([Efron et al., 2001](#); [Efron, 2004](#)) and the q -value ([Storey, 2003](#)).

2.1 Classical model-specific decision theory approaches

Using a Frequentist perspective, [Sun and Cai \(2007\)](#) frame the multiple testing problem in a compound decision theory framework. This thread of research considers controlling the marginal FDR, using the fact that, under weak conditions, $\text{mFDR} = \text{E}(\text{FDR}) + \mathcal{O}(M^{-1/2})$ ([Genovese and Wasserman, 2002](#)). [Sun and Cai \(2007\)](#) note that two approaches can be taken to address the multiple testing problem. First, one can set out with the goal of separating the non-null hypotheses from the nulls, using a weighted classification approach. In other words, the decision rule δ is constructed by minimizing the classification risk $\text{E}[L_\lambda(\theta, \delta)]$, where the loss function is

$$L_\lambda(\theta, \delta) = \frac{1}{M} \sum_{i=1}^M \left\{ \lambda(1 - \theta_i)\delta_i + \theta_i(1 - \delta_i) \right\}; \quad (2)$$

here, $\lambda > 0$ is the loss attached to a false positive error (relative to a false negative error). Alternatively, one can set out with the goal of discovering as many true findings as possible while incurring a low proportion of false positive findings: in other words, find δ with the smallest false non-discovery rate (FNR) among all rules with the FDR bounded by $\alpha \in (0, 1)$. [Sun and Cai \(2007\)](#) go on to show that these two approaches are equivalent as long as a monotone likelihood ratio condition is satisfied; that is, the optimal solution to the classification problem (where λ depends on the desired α) is also optimal for the multiple testing approach, in the sense that the classification rule yields the smallest marginal false negative rate (mFNR) among all procedures that bound $\text{mFDR} \leq \alpha$.

Unfortunately, proofs for the optimality of all of these procedures rely on the notion of independent hypotheses, and the optimality is called into question when the hypotheses are instead dependent. On one hand, [Benjamini and Yekutieli \(2001\)](#) show that FDR is controlled at the stated level for dependent hypotheses using either BH or AP. However, on the other hand, [Efron \(2007\)](#) found that non-zero correlation between tests can result in testing procedures that are either too conservative or too anti-conservative; [Schwartzman and Lin \(2011\)](#) show that the procedure can fail to be consistent as the number of tests grows under certain types of dependence. [Sun and Cai \(2009\)](#) also note that in dealing with the effects of correlation on an FDR procedure, the efficiency of the procedure should be the focus (not just the validity), and that failing to model any known dependence structure can impact the optimality of the procedure. The decision rules of [Benjamini](#)

and Hochberg (1995), Benjamini and Hochberg (2000), Efron et al. (2001), and Sun and Cai (2007) are simple, meaning that δ_i is a function only of Z_i ; i.e., $\delta_i(\mathbf{Z}) = \delta_i(Z_i)$, and therefore symmetric, meaning that $\delta(\tau(\mathbf{Z})) = \tau(\delta(\mathbf{Z}))$ for all permutation operators τ (Sun and Cai, 2007). It is easy to imagine that in the case of correlated hypotheses, compound decision rules (i.e., decision rules δ such that δ_i depends on the other $Z_j, j \neq i$) are preferred in that they might be able to identify non-nulls with a smaller signal by pooling information across tests. For example, when hypotheses are positively correlated within a temporal or spatial domain, one would expect that the non-null θ_i would appear in groups or clusters (Sun and Cai, 2009).

As a result, Sun and Cai (2009) extend the compound decision framework for multiple testing in the presence of dependence. Specifically, modeling the unknown θ_i as random effects arising from a hidden Markov model (HMM), Sun and Cai (2009) prove that the optimal classification rule for the loss function (2) is of the form $\delta_i = I(T_i < t_\lambda)$, where

$$T_i = P_{\xi}(\theta_i = 0 | \mathbf{Z}) \quad (3)$$

is the so-called “oracle statistic” and ξ is a vector of all hyperparameters in the HMM. It is important to note that the derivation of (3) in Sun and Cai (2009) as the oracle statistic is specific to the HMM framework. Furthermore, because the HMM satisfies a monotone likelihood ratio condition, T_i is also the optimal statistic for the multiple testing problem, in that $\delta_i = I(T_i < t_\lambda)$ yields the smallest mFNR subject to $\text{mFDR} \leq \alpha$. The relationship between λ and α can be seen by writing the decision rule as a step-up procedure (like BH): first, rank the oracle statistics $T_{(1)} \leq \dots \leq T_{(M)}$, and find

$$r = \max \left\{ j : \frac{1}{j} \sum_{i=1}^j T_{(i)} \leq \alpha \right\}; \quad (4)$$

then, reject $H_{(1)}, \dots, H_{(r)}$. In practice, of course, the T_i (and hence the $\{\delta_i\}$ and r) are unknown: Sun and Cai (2009) outline a data-driven procedure that uses a plug-in estimate $\hat{\xi}$ to estimate $\hat{T}_i = P_{\hat{\xi}}(\theta_i = 0 | \mathbf{Z})$ and therefore determine r by replacing $T_{(i)}$ with $\hat{T}_{(i)}$ in (4). Since the estimated oracle test statistic for the i th hypothesis depends on the entire vector of data, Sun and Cai (2009) note that the decision rule is neither simple nor symmetric.

Two recent papers by Sun et al. (2015) and Shu et al. (2015) extend the work of Sun and Cai (2009) to provide similar results for spatial random fields and multi-dimensional Markov random

fields (MRFs), respectively. In spite of the different statistical models, in both cases the oracle statistic is the same as (3) and the decision rule can be written as (4). However, model-specific proofs are required to verify that (1) the classification risk is indeed minimized by $\delta_i = I(T_i < t_\lambda)$, and (2) the optimal classification (oracle) statistic satisfies a monotone likelihood ratio condition and hence yields the smallest mFNR among all procedures with $\text{mFDR} \leq \alpha$ (here, both mFNR and mFDR are defined in a Frequentist sense). Furthermore, *estimation* of the oracle statistic T_i is, of course, model-specific. Sun and Cai (2009) use random effect prediction conditional on hyperparameter estimates: in the HMM, conditional on $\hat{\xi}$, the oracle statistic can be expressed in terms of forward and backward density variables, which can be calculated recursively. Sun et al. (2015) also conduct random effect prediction (albeit marginalizing over hyperparameters), but, since there is no longer an iterative formula for calculating the \hat{T}_i for a Gaussian random field, they instead utilize the Bayesian computational framework (i.e., Markov chain Monte Carlo) as a way to “extract information effectively from large spatial data sets” and implement their data-driven procedure.

Both Sun and Cai (2009) and Sun et al. (2015) conduct simulation studies to verify that their approach outperforms traditional FDR procedures (e.g., BH and AP) when simulated data arise from the true statistical model (i.e., HMM or Gaussian random field). However, Sun et al. (2015) also find that “the precision of [their] testing procedure shows some sensitivity to model misspecification.”

2.2 Bayesian decision theory approaches for general models

In order to move away from the model-specific procedures outlined in Section 2.1, we are motivated to consider fully Bayesian approaches to the multiple testing problem, first presented by Newton et al. (2004), Müller et al. (2004), and Müller et al. (2006). Whereas the Frequentist FDR is defined as $\text{FDR} \equiv E(\text{FDP})$ (the expectation here is taken with respect to the data, over repeated experiments), Müller et al. (2004) define a Bayesian FDR $\overline{\text{FDR}} \equiv E(\text{FDP} | \mathbf{Z}) = \int \text{FDP}(\boldsymbol{\delta}, \boldsymbol{\theta}) dp(\boldsymbol{\theta} | \mathbf{Z})$ (i.e., the posterior expected FDP), where the expectation is with respect to the posterior distribution of the unknown states conditional on the data, and we write $\text{FDP}(\boldsymbol{\delta}, \boldsymbol{\theta})$ to reiterate that FDP is a function of the true states and the decision rule. Conditioning on the data and marginalizing with respect to $\boldsymbol{\theta}$, Müller et al. (2004) show that $\overline{\text{FDR}} = [\sum_{i=1}^M \delta_i \pi_i] / [\sum_{i=1}^M \delta_i + \epsilon]$,

where $\pi_i = P(\theta_i = 0 | \mathbf{Z})$ is the posterior probability that the i th hypothesis is null and ε is a small constant to ensure that $\overline{\text{FDR}}$ is zero when $\sum_{i=1}^M \delta_i = 0$. A similar expression can be obtained for the Bayesian FNR, $\overline{\text{FNR}} = [\sum_{i=1}^M (1 - \delta_i)(1 - \pi_i)] / [M - \sum_{i=1}^M \delta_i + \varepsilon]$, as well as count versions $\overline{\text{FD}} = \sum_{i=1}^M \delta_i \pi_i$ and $\overline{\text{FN}} = \sum_{i=1}^M (1 - \delta_i)(1 - \pi_i)$.

A Bayesian decision criteria that is similar in nature to the Frequentist approaches in Section 2.1 is to minimize the $\overline{\text{FNR}}$ subject to the constraint that $\overline{\text{FDR}} \leq \alpha$. Müller et al. (2004) show that the optimal decision rule for this criteria is $\delta_i^* = I(\pi_i < t_\alpha^*)$, where the threshold depends on the desired α . Interestingly, this decision rule can be written like (4) in Sun and Cai (2009) and Sun et al. (2015): after ranking the π_i such that $\pi_{(1)} < \pi_{(2)} < \dots < \pi_{(M)}$, find

$$r_1 = \max \left\{ j : \frac{1}{j} \sum_{i=1}^j \pi_{(i)} \leq \alpha \right\}; \quad (5)$$

then $t_\alpha^* = \pi_{(r_1+1)}$, so that we reject $H_{(1)}, \dots, H_{(r_1)}$. The difference between (4) and (5) is that the former involves a probability conditional on the hyperparameters while the latter involves a probability that marginalizes over the hyperparameters. In other words, the fully Bayesian posterior probability π_i is almost the same as T_i , but accounts for uncertainty in ξ . (Note, however, that while the oracle statistic in Sun et al., 2015 is derived using a Frequentist criteria, it is calculated using a Bayesian framework and coincides exactly with (5). Their simulation study verifies that this strategy controls the Frequentist FDR.)

The optimality of (5) for controlling $\overline{\text{FDR}} \leq \alpha$ is true for “any probability model with non-zero prior probability for both the null and alternative hypotheses” (Müller et al., 2006), which is quite powerful in light of the extensive work to develop model-specific oracle procedures in the Frequentist setting (e.g., Sun and Cai, 2007, 2009; Sun et al., 2015; Shu et al., 2015). Of course, the Bayesian FDR ($\overline{\text{FDR}} \equiv E(\text{FDP} | \mathbf{Z})$) is not the same as the Frequentist FDR ($\text{FDR} \equiv E(\text{FDP})$), but Müller et al. (2004) and Müller et al. (2006) show that controlling the Bayesian FDR implies (Frequentist) FDR control when tests are independent. Unfortunately, this is not necessarily true for dependent hypotheses (Pacífico et al., 2004; Guindani et al., 2009), although the relationship between (4) and (5) suggests a similarity between the two approaches.

A benefit of the decision-theoretic framework is that classification errors can be controlled in a variety of ways, beyond just the rate of false discoveries. In addition to the decision criteria that controls the Bayesian FDR introduced in the previous paragraph by minimizing a posterior

expected loss (henceforth R_1), Müller et al. (2004) define two other decision criteria. The first (denoted R_2) is similar to the classification risk for (2):

$$R_2(\boldsymbol{\delta}, \mathbf{z}) = \lambda_1 \overline{\text{FD}} + \overline{\text{FN}} = \sum_{i=1}^M \{ \lambda_2 \delta_i \pi_i + (1 - \delta_i)(1 - \pi_i) \}.$$

This criteria minimizes the number of false negatives and false discoveries, where λ_2 represents the cost for a false discovery relative to a false negative. Like R_1 , Müller et al. (2004) show that the optimal decision rule for R_2 is a threshold rule, i.e., $\delta_i^* = I(\pi_i < t_\lambda^*)$, where the optimal threshold is $t_\lambda^* = 1/(\lambda_2 + 1)$. The second (denoted R_3) is similar in nature to R_1 , although instead of controlling the *rate* of false discoveries we control the *number* of false discoveries, i.e., R_3 minimizes the $\overline{\text{FN}}$, subject to $\overline{\text{FD}} \leq \gamma$. The optimal rule is again a threshold rule, now $\delta_i^* = I(\pi_i < t_\gamma^*)$, and we can write the optimal threshold as a step-up procedure: find

$$r_3 = \max \left\{ j : \sum_{i=1}^j \pi_{(i)} \leq \gamma \right\} \quad (6)$$

and set $t_\gamma^* = \pi_{(r_3+1)}$ so that we reject $H_{(1)}, \dots, H_{(r_3)}$. Note that by definition, R_2 and R_3 do not specifically control the false discovery rate. However, given that the optimal decision rule for both criteria is a threshold rule (like R_1), they do imply FDR control at some level determined in an indirect way via λ_2 and γ .

With all of these tools at our disposal, which should we use? On one hand, the three different decision criteria R_1 , R_2 , and R_3 allow the decision maker to choose a criteria based on their application of interest and what feels most natural. On the other hand, the criteria do not yield equivalent inference, even if the thresholds are “equivalent.” To illustrate this point, consider Figure 1, which simultaneously visualizes the three criteria by plotting artificial Bayesian posterior probabilities $\pi_i = P(\theta_i = 0 | \mathbf{Z})$ for $M = 100$ tests along with the threshold statistics corresponding to R_1 , R_2 , and R_3 . The x -axis on the bottom corresponds to the (raw) posterior probabilities π_i , and the light gray histogram in the background shows the distribution of the π_i . The x -axis along the top shows the corresponding λ_2 values, where the x -axes line up with $\pi = 1/(\lambda_2 + 1) \leftrightarrow \lambda_2 = 1/\pi - 1$. The plot shows points for three different y -axes: the left side of the plot displays the empirical cumulative distribution function (i.e., the rank/100, in blue), which is useful for illustrating R_2 ; the axes on the right show the threshold quantities for R_1 (the cumulative average of the $\pi_{(i)}$, in

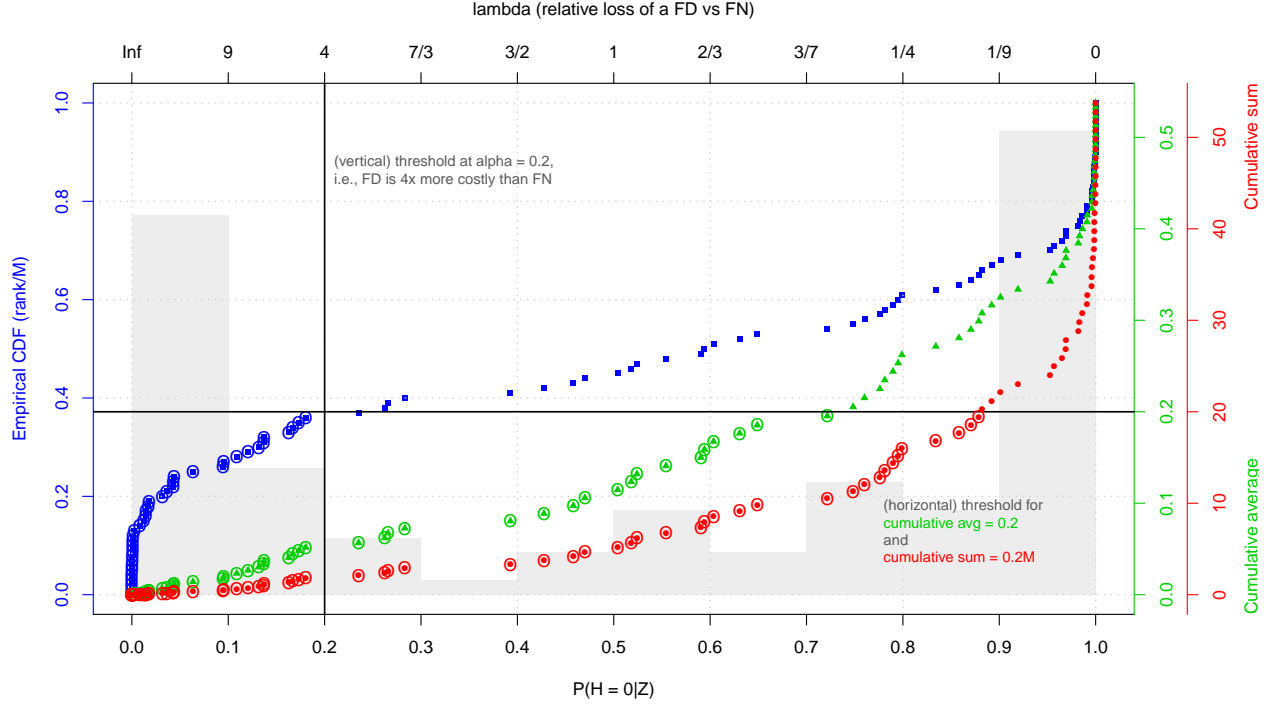


Figure 1: A comparison of the various decision criteria, for a bimodal distribution of (artificially-generated) posterior probabilities. The square points are plotted on the scale of R_2 ; the triangular points are plotted on the scale of R_1 ; the circular points are plotted on the scale of R_3 . The vertical threshold line represents the cutoff for R_2 (which thresholds the raw probabilities) when we have specified a false discovery to be 4 times more costly than a false negative. The horizontal threshold line illustrates the cutoff for both R_1 and R_3 , where we want to make sure that fewer than 20% of our discoveries are false and fewer than 20 total false discoveries, respectively.

green) and R_3 (the cumulative sum of the $\pi_{(i)}$, in red). For illustrative purposes, the plot shows the *horizontal* $\alpha = 0.2$ threshold for R_1 (meaning that we want to control $\overline{\text{FDR}}$ at 20%) as well as the “equivalent” $\gamma = 0.2 * 100 = 20$ threshold for R_3 . The *vertical* threshold is at $\lambda_2 = 1/0.2 - 1 = 4$, which indicates that controlling the $\overline{\text{FDR}}$ at 20% is equivalent to a false discovery being 4 times as costly as a false negative.

Figure 1 illustrates several key points by showing how the different threshold values from the decision criteria relate to each other. First, if we are willing to think of the R_1 and R_2 cutoffs as equivalent (i.e., that controlling $\overline{\text{FDR}}$ at 20% is equivalent to a false negative being 4 times as costly as a false discovery), then we can see that R_2 is more conservative than R_1 . This is true in general: R_2 thresholds the raw posterior probabilities π_i , while R_1 thresholds the cumulative average. Also, note that while a Bayesian model can use information from all regions to estimate

the individual posterior probabilities (see Section 3), if one uses the R_2 criteria then the distribution of the π_i is unimportant: all posterior probabilities less than the threshold are classified as rejections, regardless of how they are distributed over $(0, 1)$. Alternatively, R_1 (and R_3) considers the *cumulative* posterior probabilities when deciding the classification rule: e.g., if there are many posterior probabilities near zero, then tests with posterior probabilities much larger than α can still be rejected (in Figure 1, note that a test with $P(H = 0|\mathbf{Z}) \approx 0.72$ is rejected).

Similarly, if we are willing to think of the R_1 and R_3 cutoffs as equivalent (i.e., $\alpha = \gamma/M$), then we can see that R_1 is more conservative than R_3 (again, this is true in general). However, while equating the thresholds for R_1 and R_2 is reasonable, it is much more difficult to equate the thresholds for R_1 and R_3 ; therefore, it might not make sense to compare R_1 and R_3 . The reason for this difference is that R_1 considers a *rate* of false discoveries, while R_3 considers a *count*: as such, the total number of discoveries or rejections is very important. For example, out of 100 tests, setting out to control the $\overline{\text{FDR}}$ at 20% (using R_1) means that if 10 tests are rejected, then having 2 of those 10 *rejections* be incorrect is acceptable. This is quite different than being happy with 20 false discoveries out of 100 *tests* (which is the corresponding statement for R_3).

The bimodal distribution of posterior probabilities shown in Figure 1 illustrates that the distribution of the $\{\pi_i\}$ is important for R_1 and R_3 . Two other distributions are shown in Appendix A, comparing the decision criteria for $\{\pi_i\}$ clustered near zero (Figures A.1) and clustered near one (Figures A.2). When the π_i are clustered near zero (as in Figure 1), both R_1 and R_3 are quite aggressive and yield qualitatively similar results, rejecting tests for which the posterior probability of the null is large (i.e., tests where $\pi_i \approx 0.6$). Alternatively, when the π_i are clustered near one, R_1 is quite conservative and rejects only a few hypotheses, while R_3 is still quite liberal and rejects many hypotheses. As will be seen later, in Section 4, R_3 is always non-conservative: using this decision criteria will always result in rejecting *at least* $\lfloor \gamma \rfloor$ tests, even when all $\pi_i = 1$.

In conclusion, we reiterate that the choice of decision criteria for a specific application depends on the criteria that feels most natural for the decision maker: indeed, this is one reason that the decision-theoretic approach is so helpful. In light of the differences in R_1 , R_2 , and R_3 , our simulation study (see Section 3.2) will apply each of these decision rules and summarize the performance of each in terms of their target criteria (i.e., the realized loss, false discovery rate, and false discovery count).

3 Sensitivity to model misspecification

The classical FDR procedures in the vein of [Sun and Cai \(2007\)](#) are developed for specific data models, and unfortunately [Sun et al. \(2015\)](#) find that the optimality of the procedure is quite sensitive to model misspecification. While the Bayesian procedures of [Müller et al. \(2004\)](#) are appropriate for more general classes of models, [Newton et al. \(2004\)](#) note that the bounds on $\overline{\text{FN}}$, $\overline{\text{FD}}$, and $\overline{\text{FDR}}$ are “approximate...because [they] rest on the accuracy of the fitted model.” Furthermore, as noted in Section 2.2, the performance of the Bayesian decision rules for Frequentist FDR is not guaranteed in the presence of correlation ([Pacifico et al., 2004](#); [Guindani et al., 2009](#)). As such, we wish to understand how both model misspecification and dependence impact the performance of these decision rules for the WRAF application and, subsequently, develop a procedure that is robust to both correlation and error in specifying a statistical model for the probabilities $\{p_{ki}\}$, $k \in \{F, C\}$. First, we outline a variety of modeling frameworks for the event probabilities p_{ki} . We then conduct a simulation study to assess the performance of the various Bayesian decision rules based on R_1 , R_2 , and R_3 .

3.1 Modeling frameworks for the WRAF regions

Temperature and precipitation from the CAM5.1 climate model ensembles are aggregated monthly for each region, for both the factual and counterfactual scenarios. A climate model ensemble is a set of climate model runs such that each ensemble member has the same boundary conditions (for example, atmospheric chemistry or sea ice concentrations) but stochastically perturbed initial conditions. Denote the resulting collection of random variables $\{Y_{kil} : i = 1, \dots, M; k \in \{F, C\}; l = 1, \dots, n_{\text{ens}}\}$, where Y generically represents either average monthly temperature or total monthly precipitation and n_{ens} is the ensemble size. Formally, for an extreme event type (e.g., cold months, wet months) in region $i = 1, \dots, M$, define random variables

$$Z_{Fi} = \sum_{l=1}^{n_{\text{ens}}} I(Y_{Fil} > y_i), \quad Z_{Ci} = \sum_{l=1}^{n_{\text{ens}}} I(Y_{Cil} > y_i) \quad (7)$$

(for hot and wet extremes; replace “>” with “<” for cold and dry extremes), where the extreme event is defined in terms of a region-specific threshold y_i (e.g., exceeding a monthly average temperature threshold of 290 K; note that the event definition is common across scenarios). A natural

statistical model for these random variables is a binomial likelihood $Z_{ki} \stackrel{\text{ind}}{\sim} \text{binomial}(n_{\text{ens}}, p_{ki})$ which represents a “nonparametric” approach to estimating the event probabilities, as no assumptions need to be made regarding the behavior of the underlying climate variable (as opposed to an extreme value distribution approach). The Z_{ki} can be used to test the $H_i : RR_i \leq c$ in a variety of ways.

Classical likelihood ratio test

In order to compare new approaches with classical FDR, we first outline a method for calculating a P -value for each null hypothesis using a Frequentist likelihood ratio test. Rewriting the hypotheses in terms of the probabilities $H : p_F/p_C \leq c$ (for now ignoring the region-specific subscript), the test statistic for a likelihood ratio test considers the ratio of likelihoods for $Z_C = z_C$ and $Z_F = z_F$:

$$\lambda(z_C, z_F) = \frac{\sup_{\Theta_0} L(p_F, p_C | z_C, z_F)}{\sup_{\Theta} L(p_F, p_C | z_C, z_F)},$$

where Θ_0 is the parameter space defined by the null hypothesis and Θ is the entire (unrestricted) parameter space for p_F and p_C . The likelihood is the product of individual Binomial likelihoods:

$$L(p_F, p_C | z_C, z_F) \propto (p_F)^{z_C} (1 - p_F)^{n_C - z_C} (p_C)^{z_F} (1 - p_C)^{n_F - z_F}.$$

It can be shown that the likelihood in the denominator is maximized for the MLEs $\hat{p}_C = z_C/n_C$ and $\hat{p}_F = z_F/n_F$. Alternatively, for the numerator, the restricted MLEs are

$$(\hat{p}_C^R, \hat{p}_F^R) = \begin{cases} (\hat{p}_C, \hat{p}_F) & \text{if } \hat{p}_F > (\hat{p}_C/c) \\ (\tilde{p}_C, \tilde{p}_C/c) & \text{if } \hat{p}_F \leq (\hat{p}_C/c), \end{cases}$$

where $\tilde{p}_C = (1/4)(-b - \sqrt{b^2 - 8d})$: $b = -[c(1 + \hat{p}_F) + 1 + \hat{p}_C]$, and $d = c(\hat{p}_C + \hat{p}_F)$ (Farington and Manning, 1990). Statistical theory says $-2 \log \lambda(z_C, z_F) \xrightarrow{d} \chi_1^2$ as $n_C, n_F \rightarrow \infty$ (Θ involves two free parameters while Θ_0 has just one); thus, an asymptotic P -value is $P(\chi_1^2 > -2 \log \lambda(z_C, z_F))$. Note that when $\hat{p}_F > \hat{p}_C/c$, the likelihood ratio is 1, $-2 \log \lambda(z_C, z_F) = 0$, and the null hypothesis will never be rejected. The resulting collection of P -values can be used for an unadjusted testing procedure, classical FDR (Benjamini and Hochberg, 1995), or a Bonferroni-style family-wise error rate (FWER) procedure, which controls the probability of at least one false discovery.

Parametric Bayesian models for the risk ratio

Again using the independent binomial likelihood for the Z_{ki} in (7), the simplest Bayesian approach to modeling these probabilities is to estimate each of the p_{Fi} and p_{Ci} independently of each other and all of the other regions (an “independent across regions” model), henceforth M1. For $k \in \{F, C\}$ and $i = 1, \dots, M$, simply use a conjugate beta prior $\pi(p_{ki}) = \mathcal{B}(a_p, b_p)$ for the binomial likelihood so that the posterior is

$$\pi(p_{ki}|Z_{ki} = z_{ki}) = \mathcal{B}(z_{ki} + a_p, n_{\text{ens}} - z_{ki} + b_p); \quad (8)$$

posterior samples can be obtained by direct sampling from (8).

A more interesting Bayesian framework involves a scenario-specific hierarchical model for the probabilities

$$p_{ki} = \text{logit}^{-1}(\mu_k + \beta_{ki})$$
$$\beta_{ki} \sim G_k$$

for $k \in \{F, C\}$ and $i = 1, \dots, M$. Here, μ_k are scenario-specific (logit) means and G_k is a scenario-specific, mean-zero prior distribution for the region-specific effects. In principle, G_C and G_F need not be related; however, we model them as arising from the same parametric class but allow for different hyperparameters (and hence the subscript). Several approaches are considered:

M2 Gaussian random effects model

A random effects framework is useful for borrowing strength across the regions (“partial pooling”) without the notion of spatial dependence. As such, we can specify an exchangeable Gaussian prior $\beta_{ki} \stackrel{\text{iid}}{\sim} N(0, \tau_k^2)$. While this approach does not account for spatial dependence, each region is relatively large and therefore any dependence may be washed out by the aggregation.

M3 Skew- t random effects model

However, the Gaussian assumption may be too restrictive, in that the effects could be non-symmetric and/or heavy-tailed. Alternatively, we can use the skew- t family of distributions (Fernández and Steel, 1998; Azzalini and Capitanio, 2003; Frühwirth-Schnatter and Pyne, 2010), for which M2 is a special case. For a mean-zero random effect, this family involves three parameters: a scale parameter τ^2 , as well as ξ , which controls the degree of skewness,

and ν (the degrees of freedom), which controls the heaviness of the tails. Formally, write $\beta_{ki} \stackrel{\text{iid}}{\sim} ST(0, \tau_k^2, \xi_k, \nu_k)$. Actually, we use what [Arellano-Valle and Azzalini \(2008\)](#) call the “centered” parameterization for the skew- t distribution; see [Appendix B](#) for details.

M4 Spatially-dependent CAR model

A natural model for G_k that incorporates spatial dependence for areal data like the WRAF regions is the conditionally autoregressive (CAR) model, which models $\beta_k = (\beta_{k1}, \dots, \beta_{kM})$ as a spatial random effect (see, e.g., [Banerjee et al., 2004](#) and [Pascutto et al., 2000](#)). A CAR model relies upon the notion of a neighborhood structure; formally, for each region, define a neighborhood ∂i , which includes a subset of all regions that are specified to be “neighbors” of region i (the neighborhood definition is an important component of a CAR model; for simplicity, we take neighbors to be any regions that share a border). The benefit of a CAR model is that the joint distribution for β_k can be defined in terms of the conditional distributions $p(\beta_{ki} | \beta_k^{(-i)})$, where $\beta_k^{(-i)} = (\beta_{k1}, \dots, \beta_{k,i-1}, \beta_{k,i+1}, \dots, \beta_{kM})$. Using a Gaussian model, one representation for these conditional distributions is

$$p(\beta_{ki} | \beta_k^{(-i)}) = N \left(\frac{1}{|\partial i|} \sum_{j \in \partial i} \beta_{kj}, \frac{\tau_k^2}{|\partial i|} \right), \quad i = 1, \dots, M, \quad (9)$$

where $|\partial i| = \#$ regions in the neighborhood. (This specification is also called an *intrinsic* CAR or ICAR model.) Unfortunately, the ICAR model is an improper prior, and steps must be taken to address this problem ([Rue and Held, 2005](#); see [Appendix D](#)).

M5 Hybrid CAR/exchangeable model

The model outlined in [Leroux et al. \(2000\)](#) offers a compromise between the exchangeable random effect prior M2 and the fully spatial CAR prior M4. As outlined in [Leroux et al. \(2000\)](#), “a limitation of intrinsic autoregression is that the parameter $[\tau_k^2]$ serves both to represent overdispersion and spatial dependence”; in other words, the τ_k^2 parameter in the CAR model represents two features of the data that may be in stark contrast. A variety of strategies are used in the literature to address this problem (see, e.g., [Cressie, 1991](#); [Besag et al., 1991](#); [Leroux et al. \(2000\)](#) instead specify an approach based on “additive precisions,” in which the precision matrix of the random effects is a convex combination of the exchangeable and

CAR precision matrices:

$$\Sigma_k^{-1} \equiv \tau_k^{-2} [(1 - \lambda_k) \mathbf{I} + \lambda_k \mathbf{Q}] \quad (10)$$

where \mathbf{Q} is the CAR precision correlation matrix; $\lambda_k \in [0, 1]$ is a parameter that controls the degree of spatial dependence. Note that in this framework, $\lambda_k = 0$ corresponds to M2 while $\lambda_k = 1$ corresponds to M4; furthermore, Σ_k^{-1} is full rank for $\lambda \in [0, 1)$. Using (10), the full conditional distributions for the individual random effects are

$$p(\beta_{ki} | \beta_k^{(-i)}) = N \left(\frac{\lambda_k}{1 - \lambda_k + \lambda_k |\partial i|} \sum_{j \in \partial i} \beta_{kj}, \frac{\tau_k^2}{1 - \lambda_k + \lambda_k |\partial i|} \right), \quad i = 1, \dots, M.$$

M6 Spatial Gaussian process model

Another alternative to the CAR prior is to use a Gaussian process prior for β_k , defined for the centroids of each region (e.g., [Kelsall and Wakefield, 2002](#)). Like M5, independent random effects are a special (limiting) case of the Gaussian process prior, such that M6 can flexibly model both independent and dependent effects (unlike M4). For this approach, $\beta_k \sim N_M(\mathbf{0}, \Sigma_k)$, where

$$\Sigma_k^{ij} = \tau_k^2 \mathcal{M}_\nu \left(\frac{\|\mathbf{s}_i - \mathbf{s}_j\|}{\phi_k} \right),$$

where $\mathcal{M}_\nu(\cdot)$ is the Matérn correlation function with smoothness ν , $\mathbf{s}_i, \mathbf{s}_j$ are the three-dimensional coordinates for the centroids of regions i and j , and $\|\cdot\|$ represents Euclidean distance on \mathbb{R}^3 . In practice, since we are fitting a Gaussian process model to areal data and therefore do not observe data at very short distances, rather than trying to estimate the smoothness we instead fix it to be a constant.

Nonparametric Bayesian approach using empirically-based covariances

While the hierarchical Bayesian models outlined above are quite flexible, there could still be some degree of model-misspecification; more importantly, none of the above models can directly account for long-range dependence due to teleconnections. As an illustration of the type of correlation we might expect to see in the scenario-specific probabilities, consider [Figure 2](#), which shows the empirical correlations in the logit probability of a hot January over 1959-2014 (see below for details on how this is calculated). Note that the correlation is nearly 0.5 even at very long distances, and furthermore there appears to be a slight “bump” in the plot around 8000 km. Standard stationary

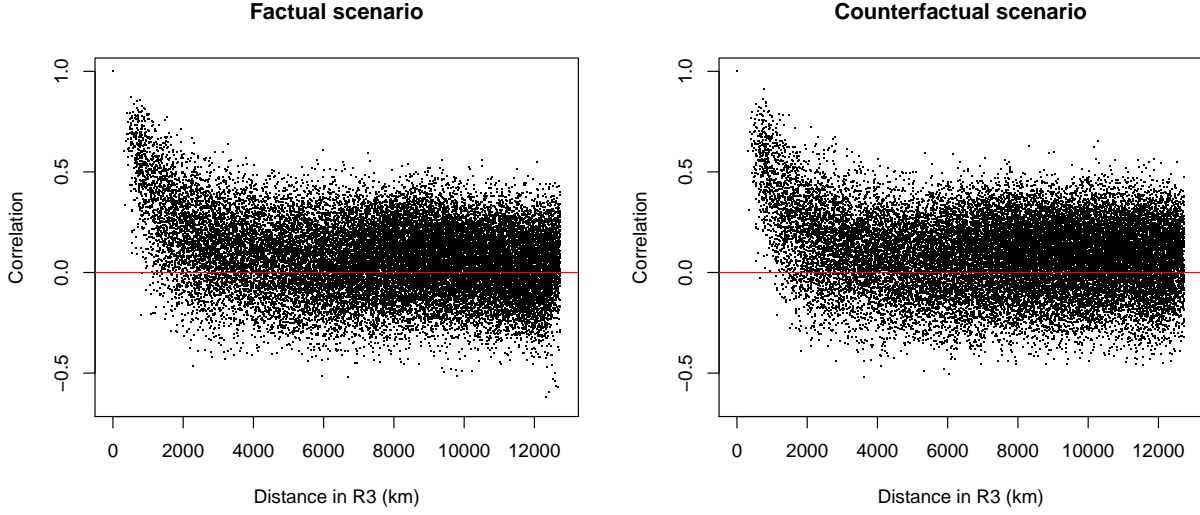


Figure 2: Empirical correlation between the logit probability of a seasonally-adjusted hot January (1959-2014; on the anomaly scale) versus distance, for both the factual (left) and the counterfactual (right) scenarios.

spatial models such as M4, M5, and M6 will likely not be able to account for these irregular dependence relationships.

Therefore, we seek a model that more robustly uses available data to estimate the covariance between the hypotheses. One way to use prior knowledge to estimate the covariances between regions is to use long time series of climate model simulations: while the WRAF is generated based on monthly simulations of the Community Atmospheric Model version 5.1 (CAM5.1; see Section 4.2 for more details), there are also historical simulations of CAM5.1 available for both climate scenarios dating back to 1959. As such, we can use the empirical relationships between the historical simulations of both the factual and the counterfactual to inform the dependence relationships among the hypotheses. Formally, we can estimate monthly probabilities $\{\hat{p}_{ki}^{(t,j)} : k \in \{F, C\}; i = 1, \dots, M; t = 1, \dots, T; j = 1, \dots, 12\}$ (t represents the year, j represents the month) using the beta-binomial model M1 (MLEs are not used because zeros are possible), where the corresponding random variables $\{z_{ki}^{(t)}\}$ are calculated using a threshold specific to each month (note: the $z_{ki}^{(t)}$ are different than the random variables introduced in (7)). Both the threshold for what is considered “extreme” and the count variables are calculated based on anomaly data (i.e., the atmospheric variables for each year are mean zero). Then, for a forecast in month j , we

have an $M \times T$ matrix of probabilities $\hat{\mathbf{p}}_k^{(j)}$ which can be used to calculate an empirical covariance on the logit scale: $\hat{\mathbf{S}}_k^{(j)} = \text{Cov} [\text{logit } \hat{\mathbf{p}}_k^{(j)}]$, where $\hat{\mathbf{p}}_k^{(j)} = \{\hat{p}_{ki}^{(t,j)} : i = 1, \dots, M; t = 1, \dots, T\}$. (Note: the correlation matrices used to create Figure 2 are from the logit $\hat{\mathbf{p}}_k^{(1)}$ for hot months.)

Unfortunately, since for our application we have $T < M$ (the historical simulations only cover $T = 56$ years and there are $M = 237$ regions), the resulting empirical estimate will not be a positive definite matrix; furthermore, it is well-known that the empirical covariance is a poor estimator for the true covariance (see, e.g., [Daniels and Kass, 2001](#); [Bickel and Levina, 2008](#)). Instead, we can use a basis function approach where the basis functions are the eigenvectors of the estimated covariance $\hat{\mathbf{S}}_k^{(j)}$, also known as empirical orthogonal functions (EOFs; see [Wikle, 2010](#) or [Cressie and Wikle, 2011](#)). EOFs are a popular strategy in modeling global climate variables, as the eigenvectors summarize the major modes of variability in a multivariate data set. Furthermore, it can be shown that the modes of variability (i.e., eigenvectors) are the same for the true covariance and a noisy estimate of the covariance (e.g., [Cressie and Wikle, 2011](#)). The main idea here is to base the current forecast on past data. While the “past” (here, 1959-2014) is not necessarily a stationary climate (especially for the factual scenario), it can be argued that the modes of variability should be approximately consistent.

As an example of the spatial patterns that we are able to capture using the EOF approach, consider the leading empirical EOFs for the logit probability of a hot January, shown in Figures A.3 and A.4 of Appendix A.

Suppressing the j notation, suppose for each month we have a set of p EOF basis functions $\mathbf{h}_{kl} = (h_{kl}(\mathbf{s}_1), \dots, h_{kl}(\mathbf{s}_M))^\top$, $l = 1, \dots, p$, for each scenario, collected into an $M \times p$ matrix $\mathbf{H}_k = (\mathbf{h}_{k1}^\top, \dots, \mathbf{h}_{kp}^\top)^\top$ (note that the EOFs are calculated separately for each scenario and event type). Then, following [Wikle \(2010\)](#), we can specify the following model for $\boldsymbol{\beta}_k = (\beta_{k1}, \dots, \beta_{kM})$:

$$\boldsymbol{\beta}_k = \mathbf{H}_k \boldsymbol{\alpha}_k + \boldsymbol{\xi}_k, \quad (11)$$

where $\boldsymbol{\alpha}$ is a p -vector of basis function coefficients with $\text{Cov}(\boldsymbol{\alpha}_k) = \boldsymbol{\Sigma}_k^\alpha$ and $\boldsymbol{\xi}_k$ is a residual vector that captures discrepancies from the EOF basis function structure. If p is large enough so that \mathbf{H}_k can account for the large-scale variability in $\boldsymbol{\beta}_k$ (later, in Section 3.2, we use $p = 30$), then we might set $\text{Cov}(\boldsymbol{\xi}_k) = \tau_k^2 \mathbf{I}_M$. Because the basis functions are orthogonal, $\boldsymbol{\Sigma}_k^\alpha$ can be diagonal. Using (11), we define two additional models, denoted M7 and M8. M7 uses a Gaussian assumption

Table 1: A summary of the models fit to each simulated data set.

Label	Model description
LRT	Classical likelihood ratio test (region-specific)
M1	Beta-binomial (independent-across-regions)
M2	Gaussian random effects
M3	Skew- t random effects
M4	Conditionally Auto-Regressive (CAR) effects
M5	Hybrid CAR/exchangeable effects
M6	Spatial Gaussian process with exponential correlation
M7	EOF-based structure, Gaussian discrepancy
M8	EOF-based structure, skew- t discrepancy

for ξ_k , i.e., $\xi_{ki} \stackrel{\text{iid}}{\sim} N(0, \tau_k^2)$ (as in M2), while M8 allows the residuals to come from a flexible skew- t distribution, i.e., $\xi_{ki} \stackrel{\text{iid}}{\sim} ST(0, \tau_k^2, \xi_k, \nu_k)$ (as in M3). Finally, we use a shrinkage prior for the EOF coefficients, i.e., $\Sigma_k^\alpha = \sigma_\alpha^2 \mathbf{I}_p$ (as in Gladish et al., 2016; see Appendix C).

A summary of all the fitted models and their labels is given in Table 1. Details on the hyperparameters and computation (via MCMC) are given in Appendices C and D.

3.2 Simulation study

To compare the performance of the decision rules introduced in Section 2 and the models introduced in Section 3.1, we perform a simulation study to explore the FDR performance for different combinations of “true” states and fitted models. For our simulation study, the number of regions will match that of the WRAF regions, i.e., $M = 237$, and we use ensemble sizes of $n_{\text{ens}} = \{50, 100, 400\}$. A total of $N_{\text{rep}} = 100$ data sets will be generated from each of six “true states” (see Table 2) that are designed to reflect the parametric and nonparametric Bayesian models; the hyperparameters for the true states will be fixed (see Appendix E), and the replicates will be drawn from the random effect distributions (as opposed to repeated binomial draws with the same effects). Complete details on the procedure for obtaining samples from each true state is provided in Appendix E; the procedure is straightforward except for generating the EOF-G and EOF-NG samples.

For each of the simulated data sets, the performance of the three decision criteria R_1 , R_2 , and R_3 will be compared by fitting each of the models outlined in Section 3.1 (see Table 1). The null

Table 2: The true states from which the simulated data sets are generated. Note: the Matérn correlation function for GP-S and GP-L has smoothness $\nu = 2$.

Label	True state
G-RE	Gaussian random effects
NG-RE	Gamma (non-Gaussian) random effects
GP-S	(Matérn) Gaussian process, short range of dependence
GP-L	(Matérn) Gaussian process, long range of dependence
EOF-G	Fixed EOF structure, Gaussian discrepancy
EOF-NG	Fixed EOF structure, gamma (non-Gaussian) discrepancy

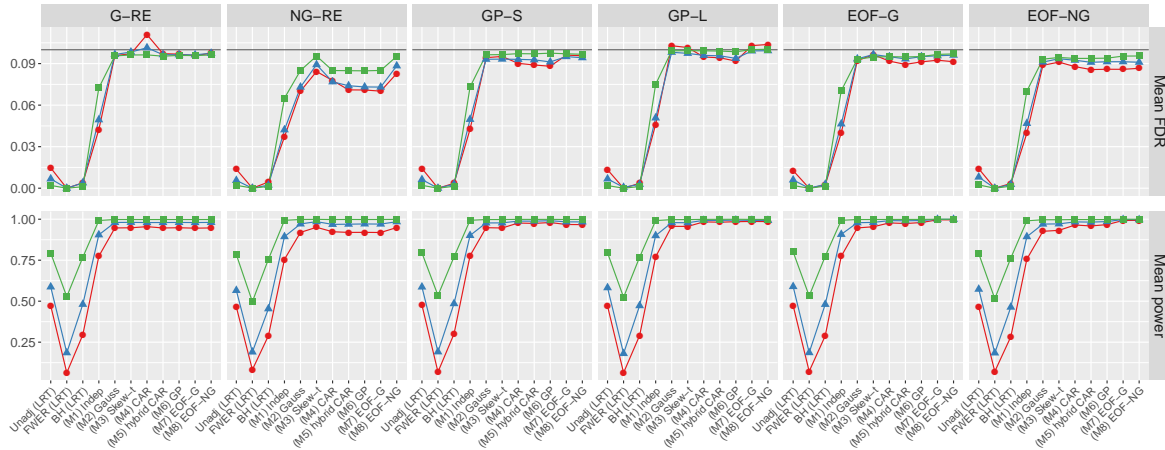
hypothesis for each simulation will use a threshold value of $c = 1$ (testing for an increase in p_F relative to p_C) while attempting to control FDR in ways comparable to the classical 0.1 significance level: for R_1 , set $\alpha = 0.10$; for R_2 , set $\lambda_2 = 1/0.1 - 1 = 9$; for R_3 , set $\gamma = 0.10M = 23.7$. Three different sets of hyperparameters will be used for each model, corresponding to cases in which most tests are true rejections (Scheme 1, ≈ 0.85), around half of tests are true rejections (Scheme 2, ≈ 0.5), and most tests are true nulls (Scheme 3, ≈ 0.15). More details on the simulations are provided in Appendix E.

Two final notes regarding the model fitting. First, for the Gaussian process model M6, note that the correlation function is set to be exponential, while the true states GP-S and GP-L have a Matérn correlation with smoothness $\nu = 2$. Second, the EOF models M7 and M8 require the initial step of estimating the EOF matrix, which is considered fixed when fitting the model. The EOFs used for true states EOF-G and EOF-NG are estimated from the historical simulations as described in Section 3.1 (the same EOFs are used for both generating data sets and model fitting). However, the benefit of the EOF framework is that it can robustly use available data to improve the model; therefore, when fitting M7 and M8 to the other true states (G-RE, NG-RE, GP-S, and GP-L), we first calculate the EOFs using $T = 56$ replicates drawn from the true state, separately for each scenario. For example, the EOFs used to fit data from GP-L would correspond to the covariance of a stationary Gaussian process with Matérn correlation function.

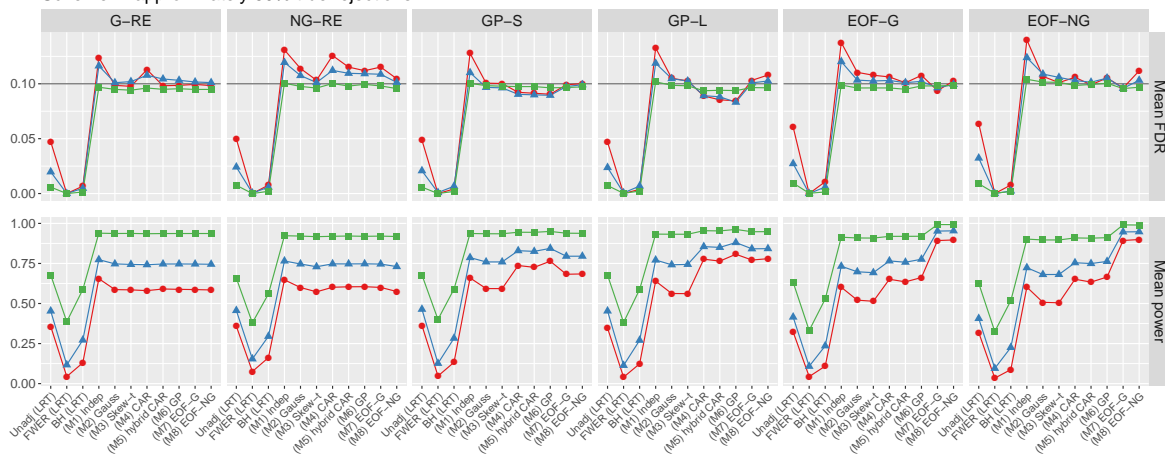
Results, summarized across simulated replicates

We present results for the R_1 criteria here, in the main text of the paper, as this decision criteria corresponds most closely with the classical notions of FDR; see Figure 3. The top, middle, and bottom sub-plots show the FDR and power (i.e., the probability of rejecting a false null) for Schemes

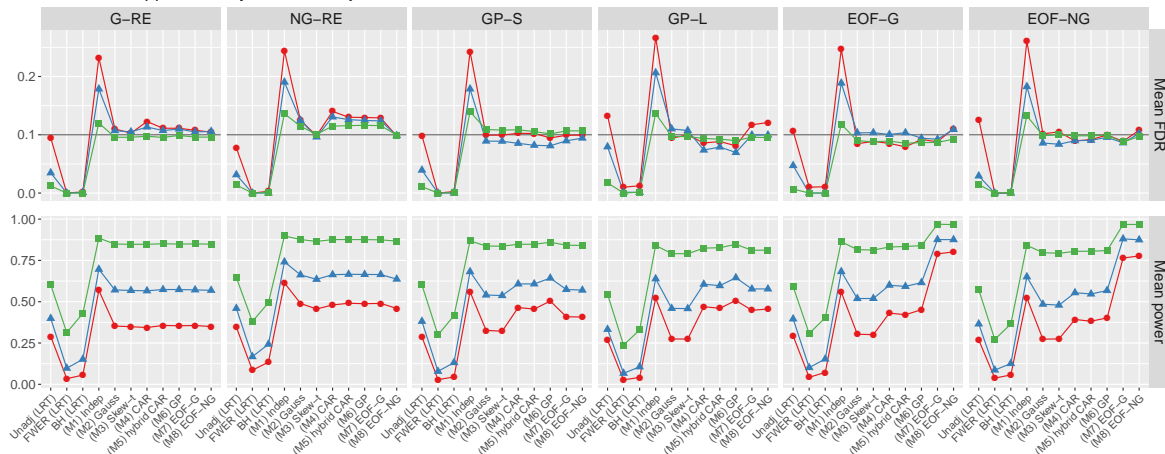
Scheme 1: approximately 85% true rejections



Scheme 2: approximately 50% true rejections



Scheme 3: approximately 15% true rejections



Ensemble size 50 100 400

Figure 3: FDR and power using the R_1 criteria, aggregated over the $N_{\text{rep}} = 100$ replicates, for schemes 1, 2, and 3. Note that the x -axis in each subgrid corresponds to the different methods/fitted models. The target of $\alpha = 0.1$ is plotted for FDR.

1, 2, and 3 (respectively), averaged over the $N_{\text{rep}} = 100$ replicated data sets. The sub-panels show the six true states, and the different methods/fitted models are shown along the x -axes.

The first observation to make is that the Frequentist P -value approaches (FWER and BH) are clearly over-conservative, such that the realized FDR is approximately zero (well below the target $\alpha = 0.1$) across all schemes, ensemble sizes, and true states. This over-conservativeness shows up in the power plots as well, with the FWER procedure in particular suffering from extremely small power, even for the largest ensemble size. Interestingly, for the Bayesian models M1-M8, each seem to do fairly well at controlling the FDR and maximizing the power (minimizing the FNR) for Scheme 1, across true states (aside from the independence model M1, which has somewhat reduced power particularly for the true states with spatial dependence).

Schemes 2 and 3 tell a different story: for each of these schemes, and across true states, the independence model is anti-conservative and fails to control the FDR (except for the largest ensemble size in Scheme 2). Otherwise, several items are noteworthy: the CAR model M4 performs poorly for the true states that do not include spatial dependence (G-RE and NG-RE); only the models that can accommodate skewness (models M3 and M8) control the FDR for the NG-RE data. Otherwise, each of the models are mostly able to control the FDR; however, major differences show up in the power. While, for example, M2 (a model without spatial dependence) is able to control the FDR for the GP-L simulations in Scheme 2, the power is significantly smaller than for a model that does accommodate dependence, e.g., M6. The EOF models M7 and M8 yield the largest power for the EOF true states by a large margin; on the other hand, these models have slightly reduced power for the GP-S and GP-L simulations (note that the reduction is larger for GP-S, which has a shorter range of spatial dependence). This is not entirely surprising since we are only using $p = 30$ EOFs to approximate a Matérn covariance, and this problem could potentially be resolved by including more EOFs.

Interestingly, the presence of spatial dependence (or lack thereof) in the simulated data has a larger effect on the power than the FDR: when the effects do not include dependence (G-RE and NG-RE), the power is roughly the same for models M2-M8. This is true even for the NG-RE effects, for which the Gaussian-based models M2, M4, M5, M6, and M7 struggle to control the FDR.

Thus, if one were to choose a “best” model for the R_1 decision criteria, the EOF model with skew- t discrepancy M8 seems to be a clear winner, performing well across schemes and true

states. This model is able to control the FDR at the nominal level for nearly every combination of true state/scheme, and almost always yields the largest power with the exception of the GP-S and GP-L true states in Schemes 2 and 3. (However, as previously mentioned, this could be addressed by including more EOFs.) In some ways, this is not surprising, since the magnitude of the EOF coefficients (together with their shrinkage prior; see Appendix C) can differentiate between cases both with and without spatial dependence and the flexibility of the skew- t residuals can capture both symmetric and non-symmetric effects. Furthermore, this approach allows us to more robustly use the data at hand (in this case, the historical CAM simulations) to capture irregular (nonstationary) spatial dependence patterns.

The story is largely the same for the R_2 and R_3 criteria (see Figures A.5 and A.6 in Appendix A): model M8 yields the smallest loss (almost always), and controls the number of FDs while minimizing the number of FNs. Therefore, we have good reason to select the Bayesian EOF-based model with skew- t discrepancy M8 combined with the decision rule of interest as the procedure that best controls the realized loss, FDR, and FD.

4 Applying the multiple testing procedure to the WRAF

Having explored the various procedures outlined in Section 2 for different combinations of true states and fitted models, we now turn to selecting a final procedure for the application at hand and applying the procedure to a real data set of climate model simulations.

4.1 Selection of decision criteria

For this application, we decided not to use the second decision criteria (R_2) because it is not clear for the WRAF what the relative loss for each type of error should be; in other words, there is no obvious way to equate the cost of a false discovery and a false negative. In deciding between R_1 and R_3 , we were initially drawn to R_3 because certain choices for the threshold γ allow us to make sure our statements are scientifically significant. In systematically conducting a set of hypotheses regarding the presence of anthropogenic influence on extreme weather, in order to conclude a significant overall (global) influence we would need to reject a null hypothesis of no anthropogenic influence for some non-zero proportion of the globe; for example, we might want

to see rejections for 5% of the globe. In other words, from a practical perspective, we might be willing to make 10 false rejections (about 5% of the 237 regions) because if we find fewer than 10 rejections then there is likely not a scientifically meaningful anthropogenic effect for the entire globe. Furthermore, in making an absolute (instead of a relative) statement about the number of false discoveries, the total number of discoveries relates to overall confidence: rejecting only a few hypotheses indicates low confidence that there is any overall anthropogenic effect, while rejecting many hypotheses indicates high confidence that there is indeed some overall anthropogenic effect.

However, upon further investigation, it became obvious that the R_3 criteria is too liberal. Returning to (6), note that this procedure will always reject *at least* $\lfloor \gamma \rfloor$ tests: in the most extreme case, where $\pi_{(i)} = 1$ for all i (meaning the null hypothesis receives all of the posterior probability), we will still have $\sum_{i=1}^{\lfloor \gamma \rfloor} \pi_{(i)} \leq \gamma$, so that in this case $r_3 = \lfloor \gamma \rfloor$. In other words, a set of tests will be rejected, even though the posterior probability that each null hypothesis is true is 1; clearly, it is quite awkward to always reject a set of hypotheses despite the evidence. Therefore, if we use the R_3 criteria, after flagging a set of null hypotheses to reject we must then determine if the results are believable. For example, if $\gamma = 10$ and we only reject ten hypotheses, then we must conclude that almost all of these are false discoveries; alternatively, if we reject 100 hypotheses, then we can be confident that most of these are true rejections. However, what if we reject 15 hypotheses? Or 20? In these “in-between” cases, we must decide when enough tests have been rejected to conclude that at least some of the rejections are true.

The R_1 criteria, on the other hand, falls more in line with traditional multiple testing procedures, in that the conclusions drawn for a set of hypotheses are more appropriately adjusted for the fact that multiple tests are being conducted. Regardless of how many tests are rejected under this criteria, we can always be sure that (in expectation) only a small proportion of these are being falsely rejected. Furthermore, while the R_3 criteria might flag some hypotheses for rejection in spite of the large posterior probability that the null is true (see the previous paragraph), the R_1 criteria will only begin flagging hypotheses for rejection if the smallest posterior probabilities of the null being true (i.e., $\pi_{(1)}$, $\pi_{(2)}$, etc.) are close to zero. A final benefit of using the R_1 criteria is that the conclusions for nested hypotheses (see Section 4.2.2) will be consistent (e.g., the procedure will only reject $H_i : RR_i^{(\text{wet})} \leq 2$ if $H_i : RR_i^{(\text{wet})} \leq 1$ is also rejected), which is not the case for R_3 .

4.2 Case study methods

Having opted to use the R_1 criteria, we set $\alpha = 0.1$ (as is done in Section 3.2). Practitioners often choose an FDR threshold based on common significance levels; here, we do the same, although there is no reason why this should be done (other than the fact that we want the FDR to be small but not too small such that we have no power). We then applied model M8 with the R_1 criteria to the WRAF for two case studies: (1) hot events in January, 2015 and (2) wet events in March, 2015. For hot events, we use a more stringent cutoff for the null hypotheses, $c_{\text{hot}} = 5$ (i.e., testing for $H_i : RR_i^{(\text{hot})} \leq 5$; this is due to the stronger anthropogenic signal for temperature), while for wet events we use $c_{\text{wet}} = 1$.

The data for estimating the p_{ki} ($k \in \{F, C\}$) consist of output from large ensembles of simulations of version 5.1 of the Community Atmospheric Model (CAM5.1) global atmosphere/land climate model, run in its conventional $\sim 1^\circ$ longitude/latitude configuration (Neale et al., 2012). Simulations have been run under the experiment protocols of the C20C+ Detection and Attribution Project (D. Stone and P. Pall, “A benchmark estimate of the effect of anthropogenic emissions on the ocean surface,” in prep.), following two historical scenarios (Angélil et al., 2017) and will be regularly updated through time as a contribution to both the C20C+ D&A project and the WRAF. The first set of simulations (for the factual scenario) is driven by observed boundary conditions of atmospheric chemistry (greenhouse gases, tropospheric and stratospheric aerosols, ozone), solar luminosity, land use/cover, and the ocean surface (temperature and ice coverage). The second set of simulations (for the counterfactual scenario) is driven by what observed boundary conditions might have been in the absence of historical anthropogenic emissions: the anthropogenic component of atmospheric chemistry is set to year-1855 values, ocean temperatures are cooled by a seasonally- and spatially-varying estimate of the warming attributable to anthropogenic emissions, and sea ice concentrations are adjusted for consistency with the ocean temperatures (Stone and Pall, in prep). Simulations within a scenario differ only in the starting conditions. The data and further details on the simulations are available at <http://portal.nersc.gov/c20c>. The simulations for both scenarios cover 01/1959 to 06/2015; the (time-varying) ensemble sizes are given in Table 3.

The event of interest for both case studies (used to define the region-specific probabilities p_{ki}) is the occurrence of a month that is more extreme than the third most extreme event expected over

Table 3: CAM5.1 ensemble sizes from January, 1959 to June, 2015 used in the case studies.

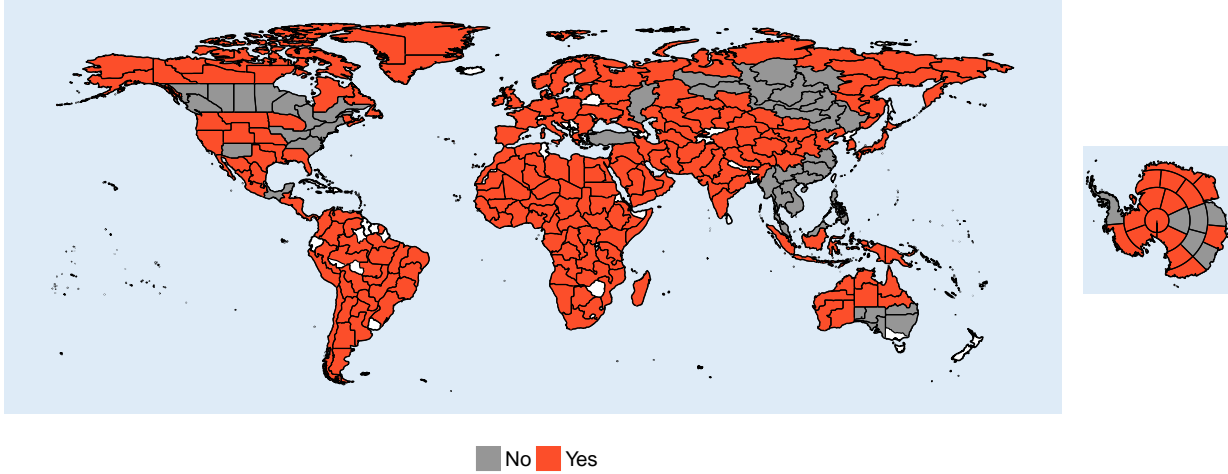
Time range	Ensemble size
01/1959 to 12/1996	50 (both factual and counterfactual)
01/1997 to 12/2009	100 (both factual and counterfactual)
01/2010 to 12/2013	400 (both factual and counterfactual)
01/2014 to 10/2014	100 (both factual and counterfactual)
10/2014 to 06/2015	98 (factual)
10/2014 to 06/2015	99 (counterfactual)

the preceding 30 year period. In other words, for a forecast in 2015, the event definition is the 1-in-10 year event which, for hot and wet months, corresponds to the 0.9 quantile of average monthly temperature or precipitation for each region in the factual simulations over 1985-2014 (specifically using the monthly measurements from the 50-member ensemble that covers this entire period). Using a moving time period of fixed length (30 years) ensures that we have accounted for climate change and that the events we consider are extreme in the “current” climate. All of the historical CAM simulations (both the factual and counterfactual) from the entire 1959-2014 period are used to calculate the EOFs following the procedure outlined in Section 3.2; the simulations from 2015 are used to fit the statistical model and classify the hypotheses. Otherwise, all prior specifications and computation via MCMC are the same as described in the simulation study. Note that an implicit assumption of this application is that the CAM5.1 simulations are suitable for evaluating changes in the probability of extremes (Angélil et al., 2016, 2017).

4.2.1 Results for a single set of hypotheses

The results for each case study are shown in Figure 4. Even with a larger cutoff for hot Januarys (testing for a five-fold increase as opposed to simply an increase), an overwhelming majority of the regions (193 of 237) have experienced a large degree of anthropogenic warming, with only a few regions in North America, Southeast Asia, central Russia, and southeast Australia failing to provide conclusive evidence of a five-fold increase in occurrence probability of a hot month (every region has conclusive evidence against the null hypothesis when the cutoff is relaxed to $c_{\text{hot}} = 1$). The results are more varied for wet events in March, as there are many regions with and without conclusive evidence against the null hypothesis. Many regions in the northern extratropics (mid to

Conclusive evidence for 5X increase in probability of a hot January in 2015



Conclusive evidence for an increase in probability of a wet March in 2015

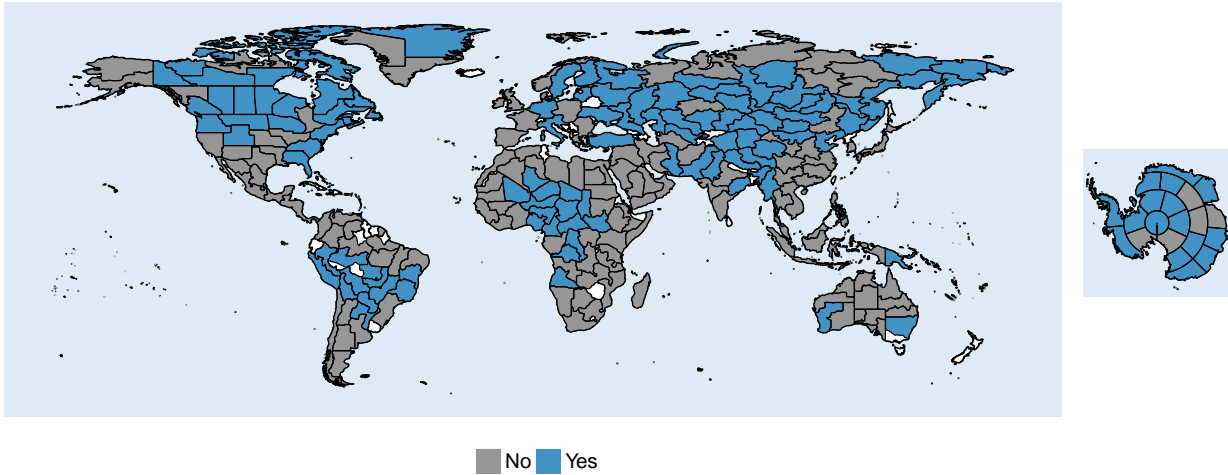


Figure 4: Results of testing a collection of hypotheses $H_i : RR_i^{(\text{hot})} \leq 5$ (top; i.e., determining if there is conclusive evidence for a five-fold increase in the probability that January, 2015 will have an average temperature that exceeds the third hottest expected January over 1985-2014) and $H_i : RR_i^{(\text{wet})} \leq 1$ (bottom; i.e., determining if there is conclusive evidence for an increase in the probability that the total precipitation in March, 2015 will exceed the third wettest expected March over 1985-2014). The white areas (e.g. New Zealand, Zimbabwe) do not satisfy criteria for fitting into political regions of the target 400,000-900,000 km² range as described in Stone (in prep) and are not analyzed here.

high latitudes) have an increased probability of a wet event as a result of anthropogenic emissions. An increased probability is the general tendency along an equatorial band as well (although not in Southeast Asia), while the subtropics (arid regions in the tropics) generally lack conclusive

evidence; a notable exception is over the northern Sahel, which may indicate an earlier advance of the West African monsoon in this climate model due to anthropogenic emissions (Lawal et al., 2016).

4.2.2 Capturing the existence and magnitude of anthropogenic influence

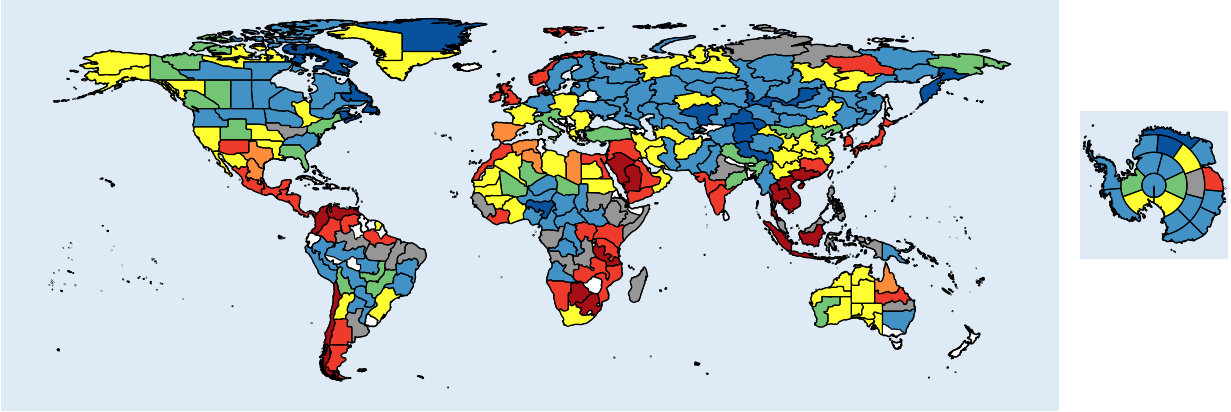
Both the current and planned upcoming versions of the WRAF actually conduct more than one set of hypothesis tests for each forecast: several different thresholds are used (e.g., $c_{\text{wet}} = 1$ versus $c_{\text{wet}} = 2$) in conjunction with several different types of null hypotheses (e.g., $H_i : RR_i^{(\text{wet})} \leq c_{\text{wet}}$ versus $H_i : RR_i^{(\text{wet})} \geq c_{\text{wet}}$). The purpose of these categories is to make statements that combine confidence in the change in probability as well as the magnitude of this change. As such, the forecast actually involves “multiple-multiple testing,” in that we now have multiple sets of M hypotheses to test. This can be accomplished in our framework by simply conducting the classification procedure several times; recall from Section 4.1 that using R_2 yields consistent results for nested hypotheses (unlike R_3). The testing adjustment is done separately for each category, and therefore the existence of any possible false discoveries can be interpreted within each category.

As an example of what the attribution forecast looks like for multiple categories, see Figure 5. A benefit of the Bayesian framework is that we can first test for the *absence* of an anthropogenic effect using a null hypothesis like $H_i^{(1)} : RR_i^{(\text{wet})} \leq l_{\text{absence}} \cup RR_i^{(\text{wet})} \geq u_{\text{absence}}$, where u_{absence} and l_{absence} are upper and lower limits, respectively, for an interval including 1 that defines “no anthropogenic influence.” Regions where we can reject $H_i^{(1)}$ display strong evidence that anthropogenic forcings have not changed the probability of extreme precipitation. Otherwise, the other null hypotheses of interest are

$$H_i^{(2)} : RR_i^{(\text{wet})} \geq 1/2; \quad H_i^{(3)} : RR_i^{(\text{wet})} \geq 1; \quad H_i^{(4)} : RR_i^{(\text{wet})} \leq 1; \quad H_i^{(5)} : RR_i^{(\text{wet})} \leq 2;$$

being able to reject these hypotheses indicates conclusive evidence that the probability of extreme precipitation is decreased by a factor of two, decreased, increased, or increased by a factor of two (respectively). There is clearly some overlap between $H_i^{(1)}$ and both $H_i^{(3)}$ and $H_i^{(4)}$; to reflect this, we create two additional categories to indicate regions that reject both $H_i^{(1)}$ and $H_i^{(3)}$ (orange, indicating that while there is most likely no change in the probability there is some evidence for a decrease) as well as both $H_i^{(1)}$ and $H_i^{(4)}$ (green, indicating that while there is most likely no

Conclusive evidence for changes in probability of a wet March in 2015



Conclusive evidence for changes in probability of a wet March in 2015

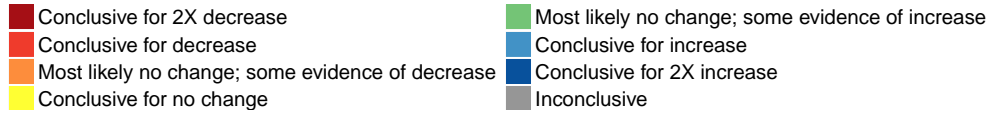
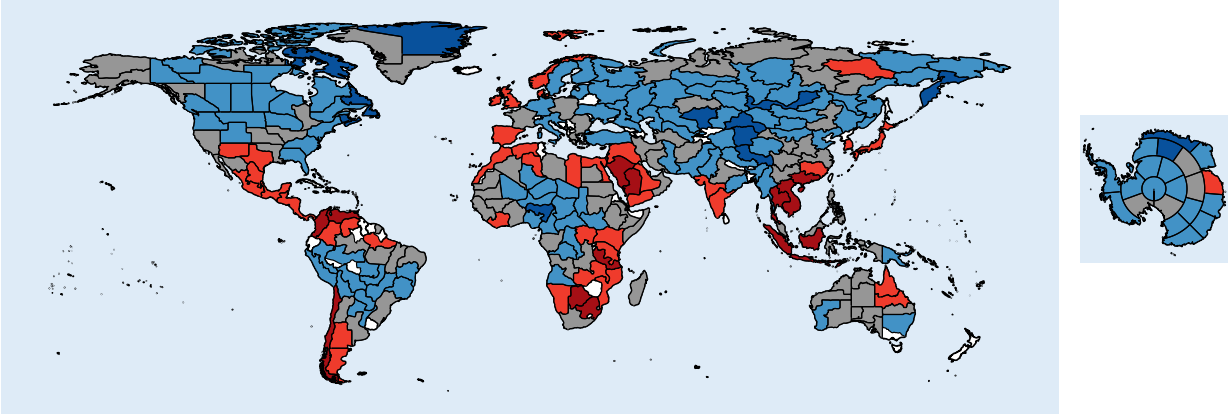


Figure 5: Results of testing multiple hypotheses per region, in order to capture the magnitude and direction of the effect of anthropogenic influence. Top: $1/2 \leq RR_i \leq 2$ defines the “Conclusive for no change” category; bottom: $2/3 \leq RR_i \leq 3/2$ defines the “Conclusive for no change” category.

change in the probability there is some evidence for an increase). A final category (shown in gray) identifies regions that fail to reject any of the hypotheses and are thus classified as inconclusive.

Maps of these multi-category results are shown in Figure 5, where we use both a wide interval $l_{\text{absence}} = 1/2$ and $u_{\text{absence}} = 2$ as well as narrower limits $l_{\text{absence}} = 2/3$ and $u_{\text{absence}} = 3/2$. Angéilil et al. (2017) provide justification for using the narrower interval $(2/3, 3/2)$ as the definition of “no anthropogenic influence”; however, the somewhat limited ensemble sizes (≈ 100) in this case

study prevent us from conclusively finding no change for the narrower interval (bottom, Figure 5). The wider interval used for the top panel of Figure 5, on the other hand, concludes that a fairly large proportion of the map experiences no change. Clearly, being able to conclude that extreme probabilities are unchanged between the two climate scenarios depends heavily upon both the ensemble size and width of the interval that defines “no change,” and this tradeoff can result in significant qualitative differences. Given that the geometric range spanned by $(2/3, 3/2)$ is about half that spanned by $(1/2, 2)$, a quadrupling of the ensemble size (i.e., increased to 400 members) would be expected to result in the identification of a substantial number of regions in the “no change” categories, in contrast to when only 100 members are available (as also noted in analysis of 2013 events when the ensemble sizes are larger; not shown).

5 Conclusions

In this paper, we have developed a hierarchical Bayesian modeling framework for estimating the probability of extreme events and the risk ratio over a large collection of land-regions, as well as a decision theoretic procedure that allows us to flexibly control the number of false discoveries while maximizing the number of true discoveries. The Bayesian hierarchical model robustly uses historical climate model simulations to estimate irregular (nonstationary) dependence patterns among the hypotheses and can account for non-Gaussian behavior in the region-specific effects. Furthermore, we show that the modeling framework maintains false discovery control even when the true data-generating mechanism arises from a completely different class of statistical models. Finally, we apply our robust statistical model to a real data set used for making seasonal forecasts for the Weather Risk Attribution Forecast. Moving forward, we plan to operationalize our procedure as described in Section 4.2 to replace the current ad hoc presentation of the forecast.

We have demonstrated the application of our procedure across regions and with multiple hypotheses for each region. However, we have not applied it across event types (e.g. hot, cold, wet, and dry events for a single region) or across multiple months. Application across regions makes sense for several reasons. First, events are presented in global maps of these regions (as in Figure 5) and thus are not only providing information on each region individually but also on the aggregate of all of the regions. Second, even with some correlation across the regions, there remains a large “effective sample size” of tests, whereas testing across event types would yield a small num-

ber of tests (such that a multiple testing adjustment has less value). While testing across multiple months (e.g., all months or only the same calendar month from a given period of years) may provide a moderate number of tests, it would be hard to fit into the monthly operational design of the WRAF. Continual updating of past calculations, as further months become available, would pose a presentation and communication challenge. However, in a more retrospective research framework, studying events over a decade for instance, testing across months as well as, or instead of, regions could make sense.

Finally, while the final version of the forecast will use the $M = 237$ regions shown in Figure 4 (where each region is approximately 0.5 million km²; these are the WRAF05 regions), the WRAF will also provide a forecast for aggregates of these regions: 68 regions comprising 2 million km² each (WRAF2); 30 regions comprising 5 million km² each (WRAF5); and 12 regions comprising 10 million km² each (WRAF10). As a demonstration of how our procedure will perform for a smaller number of regions, we conducted a simulation study similar to the one in Section 3.2 using the larger WRAF2 regions ($M = 68$; these regions are slightly modified from the current version of the WRAF, which has 58 regions); these results are shown in Appendix A.2. Results for the WRAF2 regions are approximately consistent with the simulation study results for the smaller, WRAF05 regions.

Acknowledgements

This research was supported by the Director, Office of Science, Office of Biological and Environmental Research of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231 as part of their Regional & Global Climate Modeling Program (RGCM) and used resources of the National Energy Research Scientific Computing Center (NERSC), also supported by the Office of Science of the U.S. Department of Energy, under Contract No. DE-AC02-05CH11231.

References

- Allen, M. (2003). Liability for climate change. *Nature*, 421:891–892.
- Angélil, O., Perkins-Kirkpatrick, S., Alexander, L. V., Stone, D., Donat, M. G., Wehner, M., Shiogama, H., Ciavarella, A., and Christidis, N. (2016). Comparing regional precipitation and temperature extremes in climate model and reanalysis products. *Weather and Climate Extremes*, 13:35–43.
- Angélil, O., Stone, D., Wehner, M., Paciorek, C. J., Krishnan, H., and Collins, W. (2017). An

- independent assessment of anthropogenic attribution statements for recent extreme temperature and rainfall events. *Journal of Climate*, 30(1):5–16.
- Angéilil, O., Stone, D. A., Tadross, M., Tummon, F., Wehner, M., and Knutti, R. (2014). Attribution of extreme weather to anthropogenic greenhouse gas emissions: sensitivity to spatial and temporal scales. *Geophys. Res. Lett.*, 41:2150–2155.
- Arellano-Valle, R. B. and Azzalini, A. (2008). The centred parametrization for the multivariate skew-normal distribution. *Journal of Multivariate Analysis*, 99(7):1362 – 1382. Special Issue: Multivariate Distributions, Inference and Applications in Memory of Norman L. Johnson.
- Arent, D. J., Tol, R. S. J., Faust, E., Hella, J. P., Kumar, S., Strzepek, K. M., Tóth, F. L., Yan, D., and et alii (2014). Key economic sectors and services. In Field, C. B., Barros, V. R., and et alii, editors, *Climate Change 2014: Impacts, Adaptation, and Vulnerability. Part A: Global and Sectoral Aspects. Contribution of Working Group II to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, pages 659–708. Cambridge University Press.
- Azzalini, A. and Capitanio, A. (2003). Distributions generated by perturbation of symmetry with emphasis on a multivariate skew t-distribution. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(2):367–389.
- Banerjee, S., Carlin, B. P., and Gelfand, A. E. (2004). *Hierarchical Modeling and Analysis for Spatial Data*. Monographs on Statistics and Applied Probability. Chapman & Hall, London.
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):pp. 289–300.
- Benjamini, Y. and Hochberg, Y. (2000). On the adaptive control of the false discovery rate in multiple testing with independent statistics. *Journal of Educational and Behavioral Statistics*, 25(1):60–83.
- Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *The Annals of Statistics*, 29(4):1165–1188.
- Besag, J., York, J., and Mollié, A. (1991). Bayesian image restoration, with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics*, 43(1):1–20.
- Bickel, P. J. and Levina, E. (2008). Covariance regularization by thresholding. *The Annals of Statistics*, pages 2577–2604.
- Cressie, N. and Wikle, C. K. (2011). *Statistics for Spatio-Temporal Data*. Probability and Statistics. John Wiley & Sons, New York.
- Cressie, N. A. C. (1991). *Statistics for spatial data*. John Wiley & Sons, New York.
- Daniels, M. J. and Kass, R. E. (2001). Shrinkage estimators for covariance matrices. *Biometrics*, 57(4):1173–1184.

- Efron, B. (2004). Local false discovery rate. *Unpublished technical report*.
- Efron, B. (2007). Correlation and large-scale simultaneous significance testing. *Journal of the American Statistical Association*, 102(477):93–103.
- Efron, B., Tibshirani, R., Storey, J. D., and Tusher, V. (2001). Empirical bayes analysis of a microarray experiment. *Journal of the American Statistical Association*, 96(456):1151–1160.
- Farrington, C. P. and Manning, G. (1990). Test statistics and sample size formulae for comparative binomial trials with null hypothesis of non-zero risk difference or non-unity relative risk. *Statistics in Medicine*, 9(12):1447–1454.
- Fernández, C. and Steel, M. F. (1998). On Bayesian modeling of fat tails and skewness. *Journal of the American Statistical Association*, 93(441):359–371.
- Frühwirth-Schnatter, S. and Pyne, S. (2010). Bayesian inference for finite mixtures of univariate and multivariate skew-normal and skew-t distributions. *Biostatistics*, 11(2):317–336.
- Genovese, C. and Wasserman, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):499–517.
- Gladish, D. W., Kuhnert, P. M., Pagendam, D. E., Wikle, C. K., Bartley, R., Searle, R. D., Ellis, R. J., Dougall, C., Turner, R. D. R., Lewis, S. E., Bainbridge, Z. T., and Brodie, J. E. (2016). Spatio-temporal assimilation of modelled catchment loads with monitoring data in the great barrier reef. *Ann. Appl. Stat.*, 10(3):1590–1618.
- Guindani, M., Miller, P., and Zhang, S. (2009). A bayesian discovery procedure. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):905–925.
- Hansen, G., Auffhammer, M., and Solow, A. R. (2014). On the attribution of a single event to climate change. *J. Climate*, 27:8297–8301.
- Kelsall, J. and Wakefield, J. (2002). Modeling spatial variation in disease risk. *Journal of the American Statistical Association*, 97(459):692–701.
- Lawal, K. A., Abatan, A. A., Angélil, O., Olaniyan, E., Olusoji, V. H., Oguntunde, P. G., Lamptey, B., Abiodun, B. J., Shiogama, H., Wehner, M. F., and Stone, D. A. (2016). The late onset of the 2015 wet season in Nigeria. *Bull. Amer. Meteor. Soc.*, 97:S63–S69.
- Leroux, B. G., Lei, X., and Breslow, N. (2000). *Estimation of Disease Rates in Small Areas: A new Mixed Model for Spatial Dependence*, pages 179–191. Springer New York, New York, NY.
- Livezey, R. E. and Chen, W. Y. (1983). Statistical field significance and its determination by monte carlo techniques. *Monthly Weather Review*, 111(1):46–59.
- Müller, P., Parmigiani, G., and Rice, K. (2006). Fdr and bayesian multiple comparisons rules. *Johns Hopkins University, Dept. of Biostatistics Working Papers*, 115.

- Müller, P., Parmigiani, G., Robert, C., and Rousseau, J. (2004). Optimal sample size for multiple testing: The case of gene expression microarrays. *Journal of the American Statistical Association*, 99(468):990–1001.
- National Academies of Sciences, E. and Medicine (2016). *Attribution of Extreme Weather Events in the Context of Climate Change*. The National Academies Press, Washington, DC.
- Neale, R. B., Chen, C.-C., Gettelman, A., Lauritzen, P. H., Park, S., Williamson, D. L., Conley, A. J., Garcia, R., Kinnison, D. Lamarque, J.-F., Marsh, D., Mills, M., Smith, A. K., Tilmes, S. Vitt, F., Morrison, H., Cameron-Smith, P., Collins, W. D., Iacono, M. J., Easter, R. C., Ghan, S. J., Liu, X., Rasch, P. J., and Taylor, M. A. (2012). Description of the NCAR community atmosphere model (CAM 5.0). Technical report, NCAR Technical Note NCAR/TN-486+STR.
- Newton, M. A., Noueiry, A., Sarkar, D., and Ahlquist, P. (2004). Detecting differential gene expression with a semiparametric hierarchical mixture method. *Biostatistics*, 5(2):155–176.
- NIMBLE Development Team (2016). NIMBLE: An R package for programming with BUGS models, version 0.6-3.
- Pacifico, M. P., Genovese, C., Verdinelli, I., and Wasserman, L. (2004). False discovery control for random fields. *Journal of the American Statistical Association*, 99(468):1002–1014.
- Pall, P., Aina, T., Stone, D. A., Stott, P. A., Nozawa, T., Hilberts, A. G. J., Lohmann, D., and Allen, M. R. (2011). Anthropogenic greenhouse gas contribution to flood risk in England and Wales in Autumn 2000. *Nature*, 470:382–385.
- Pascutto, C., Wakefield, J., Best, N., Richardson, S., Bernardinelli, L., Staines, A., and Elliott, P. (2000). Statistical issues in the analysis of disease mapping data. *Statistics in Medicine*, 19(17-18):2493–2519.
- Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*, volume 104 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London.
- Schwartzman, A. and Lin, X. (2011). The effect of correlation in false discovery rate estimation. *Biometrika*, 98(1):199–214.
- Shu, H., Nan, B., and Koeppe, R. (2015). Multiple testing for neuroimaging via hidden Markov random field. *Biometrics*, 71(3):741–750.
- Smith, K. R., Woodward, A., Campbell-Lendrum, D., Chadee, D. D., Honda, Y., Liu, Q., Olwoch, J. M., Revich, B., Sauerborn, R., and et alii (2014). Human health: impacts, adaptation, and co-benefits. In Field, C. B., Barros, V. R., and et alii, editors, *Climate Change 2014: Impacts, Adaptation, and Vulnerability. Part A: Global and Sectoral Aspects. Contribution of Working Group II to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, pages 709–754. Cambridge University Press.
- Stone, D. A. and Allen, M. R. (2005). The end-to-end attribution problem: from emissions to impacts. *Clim. Change*, 71:303–318.

- Storey, J. D. (2003). The positive false discovery rate: a Bayesian interpretation and the q-value. *Ann. Statist.*, 31(6):2013–2035.
- Stott, P. A., Allen, M., Christidis, N., Dole, R. M., Hoerling, M., Huntingford, C., Pall, P., Perlwitz, J., and Stone, D. (2013). *Attribution of Weather and Climate-Related Events*, pages 307–337. Springer Netherlands, Dordrecht.
- Stott, P. A., Stone, D. A., and Allen, M. R. (2004). Human contribution to the European heatwave of 2003. *Nature*, 432:610–614.
- Sun, W. and Cai, T. T. (2007). Oracle and adaptive compound decision rules for false discovery rate control. *Journal of the American Statistical Association*, 102(479):901–912.
- Sun, W. and Cai, T. T. (2009). Large-scale multiple testing under dependence. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(2):393–424.
- Sun, W., Reich, B. J., Tony Cai, T., Guindani, M., and Schwartzman, A. (2015). False discovery control in large-scale spatial multiple testing. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(1):59–83.
- Ventura, V., Paciorek, C. J., and Risbey, J. S. (2004). Controlling the proportion of falsely rejected hypotheses when conducting multiple tests with climatological data. *Journal of Climate*, 17(22):4343–4356.
- Wikle, C. K. (2010). Low-rank representations for spatial processes. In Gelfand, A. E., Diggle, P. J., Fuentes, M., and Guttorp, P., editors, *Handbook of Spatial Statistics*, chapter 8, pages 107–118. Chapman and Hall/CRC, London.
- Wilks, D. S. (2016). the stippling shows statistically significant grid points: How research results are routinely overstated and overinterpreted, and what to do about it. *Bulletin of the American Meteorological Society*, 97(12):2263–2273.

A Supplemental figures

Supplemental figures for the main text are shown in Figures A.1, A.2, A.3, A.4, A.5, and A.6. Results for the simulation study with the larger WRAF regions (WRAF2 with $M = 68$) are shown in Figures A.7, A.8, and A.9.

A.1 Main text

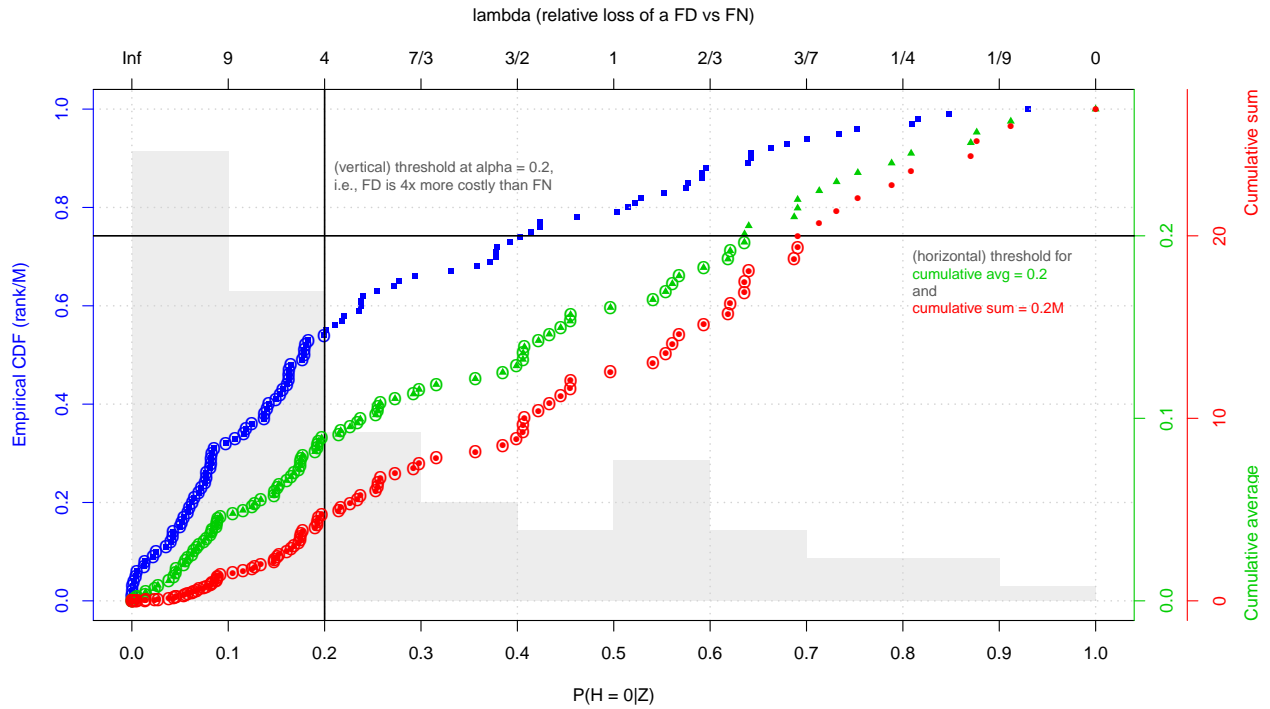


Figure A.1: A comparison of the various decision criteria, for artificially-generated posterior probabilities clustered around zero. The square points are plotted on the scale of R_2 ; the triangular points are plotted on the scale of R_1 ; the circular points are plotted on the scale of R_3 . The vertical threshold line represents the cutoff for R_2 (which thresholds the raw probabilities) when we have specified a false discovery to be 4 times more costly than a false negative. The horizontal threshold line illustrates the cutoff for both R_1 and R_3 , where we want to make sure that fewer than 20% of our discoveries are false and fewer than 20 total false discoveries, respectively.

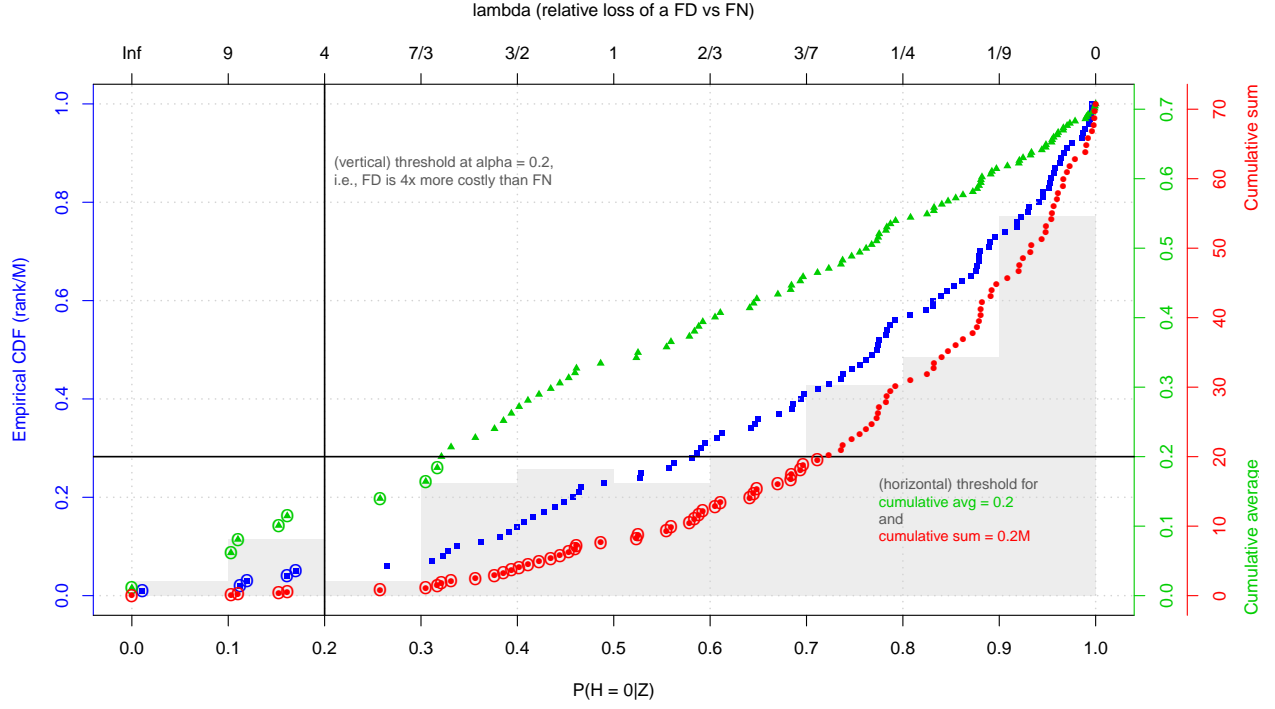


Figure A.2: A comparison of the various decision criteria, for artificially-generated posterior probabilities clustered around one. The square points are plotted on the scale of R_2 ; the triangular points are plotted on the scale of R_1 ; the circular points are plotted on the scale of R_3 . The vertical threshold line represents the cutoff for R_2 (which thresholds the raw probabilities) when we have specified a false discovery to be 4 times more costly than a false negative. The horizontal threshold line illustrates the cutoff for both R_1 and R_3 , where we want to make sure that fewer than 20% of our discoveries are false and fewer than 20 total false discoveries, respectively.

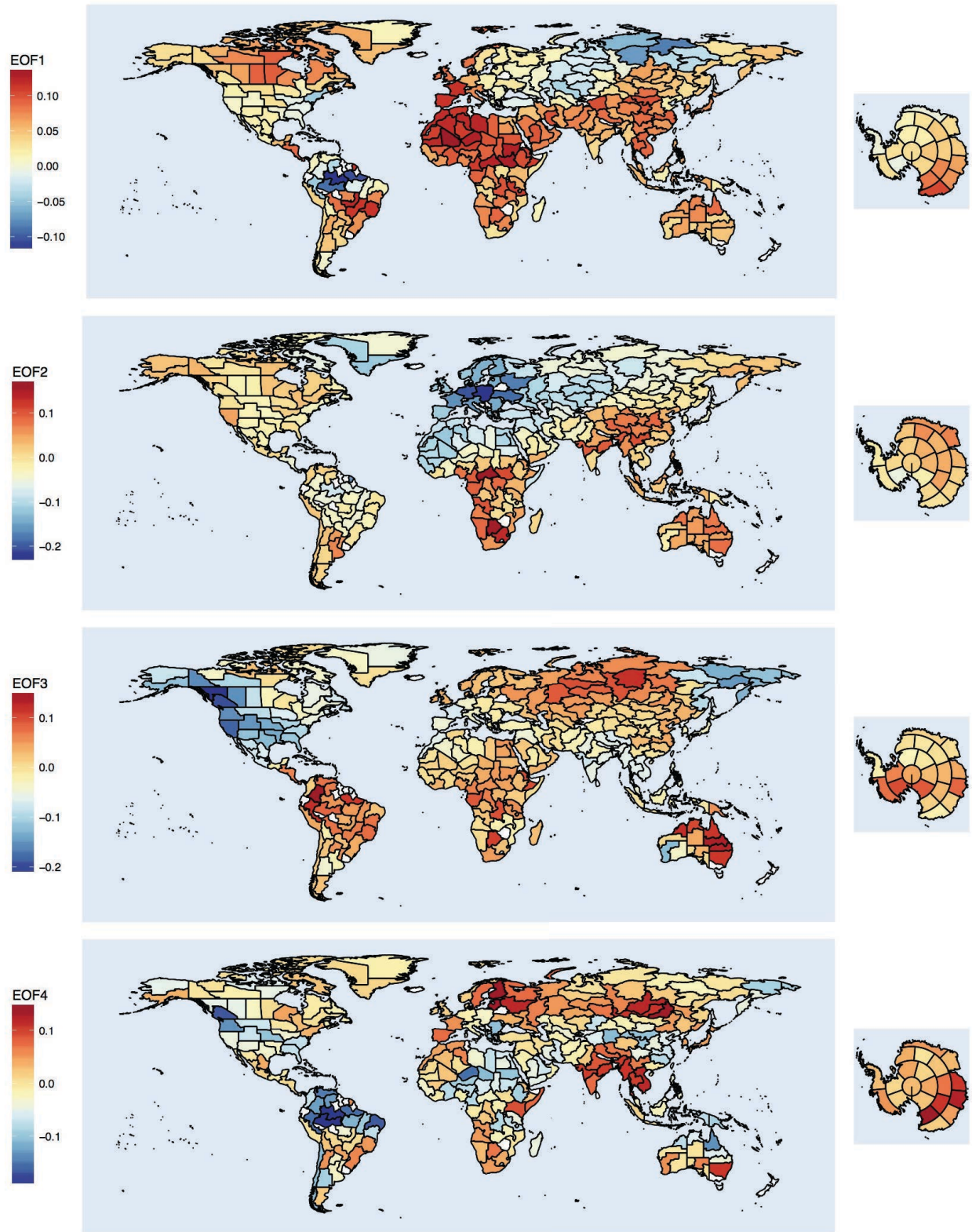


Figure A.3: The first four EOFs for the logit probability of a hot January over 1959-2014, for the factual scenario.

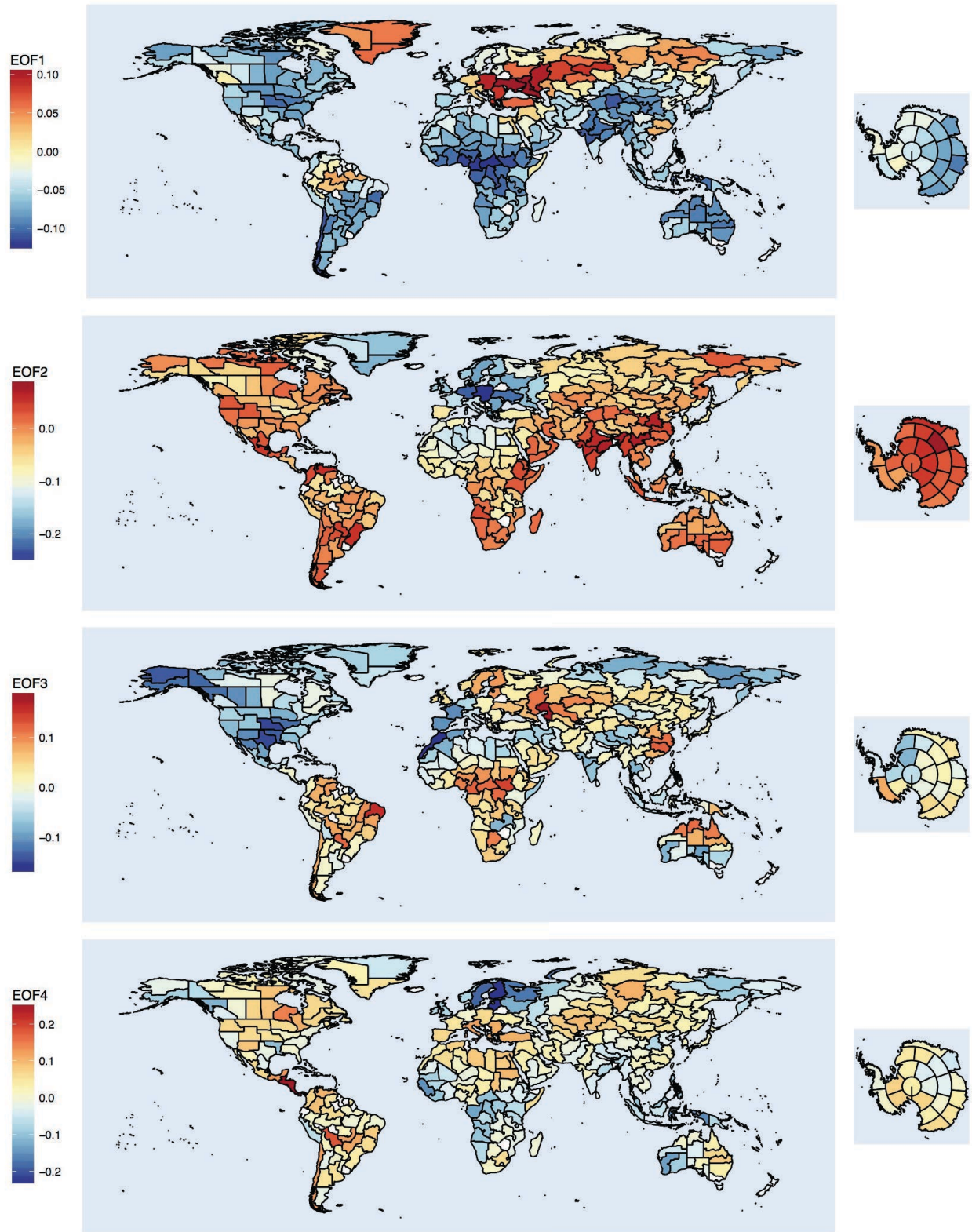


Figure A.4: The first four EOFs for the logit probability of a hot January over 1959-2014, for the counterfactual scenario.

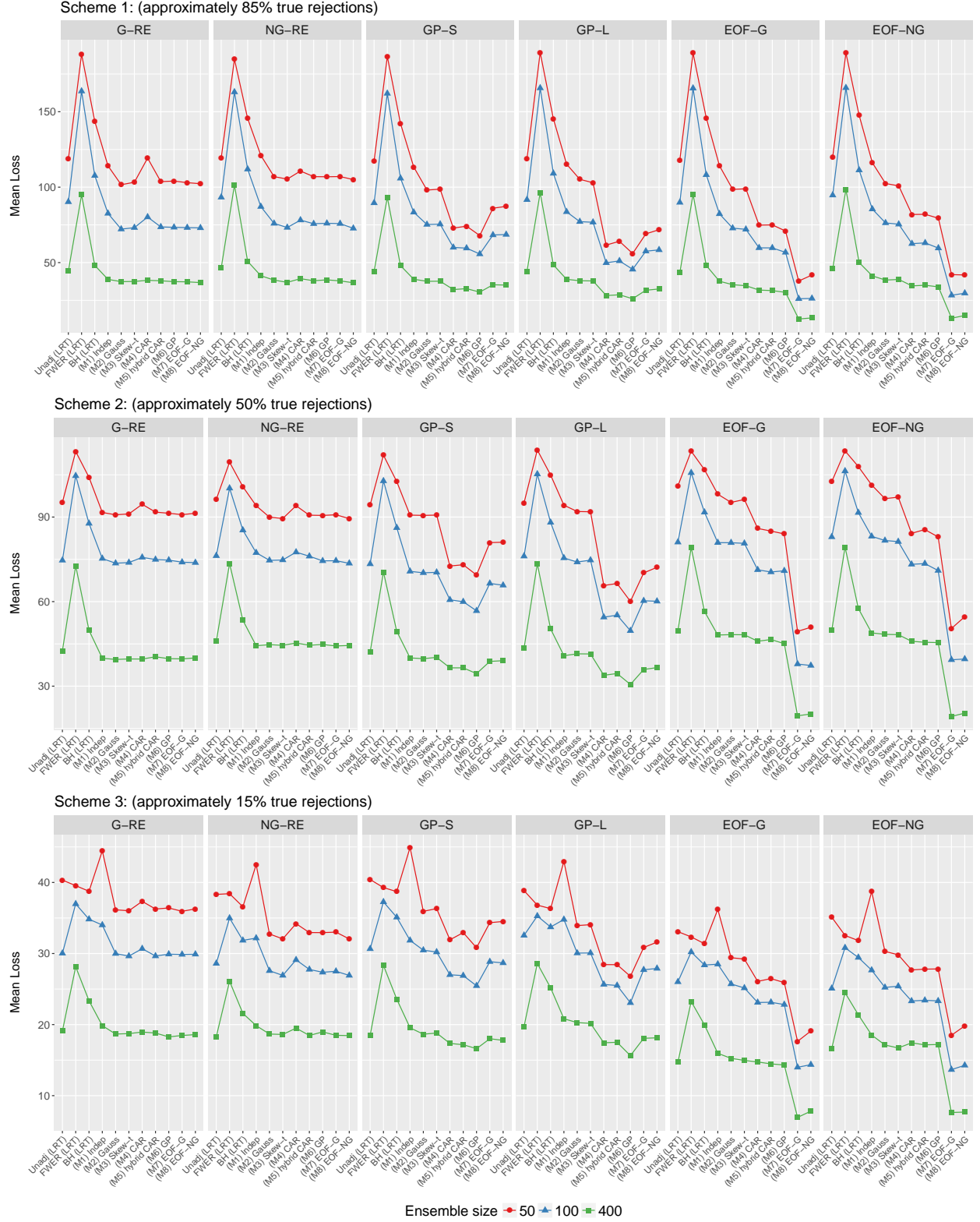
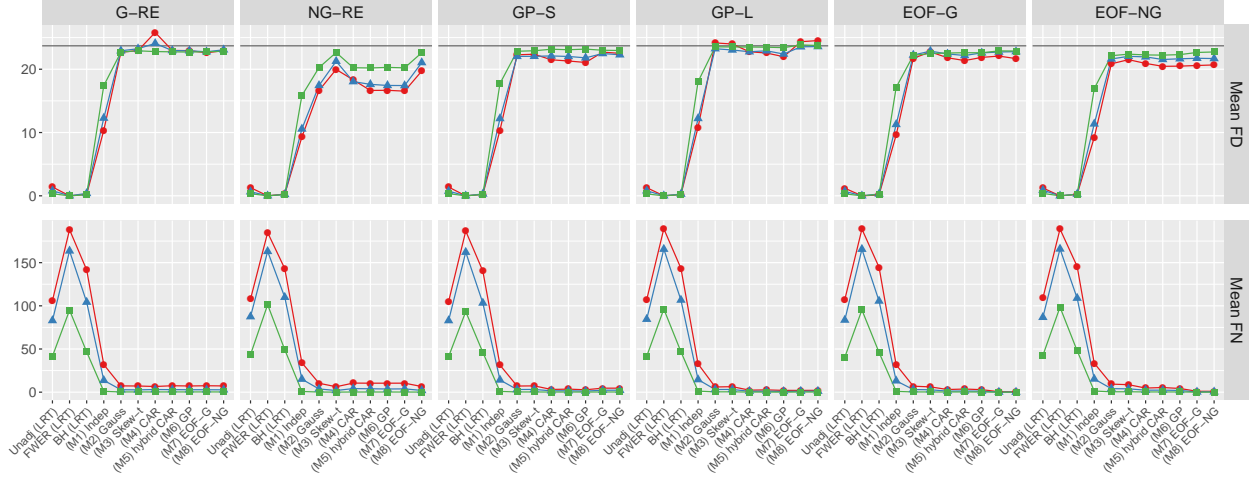
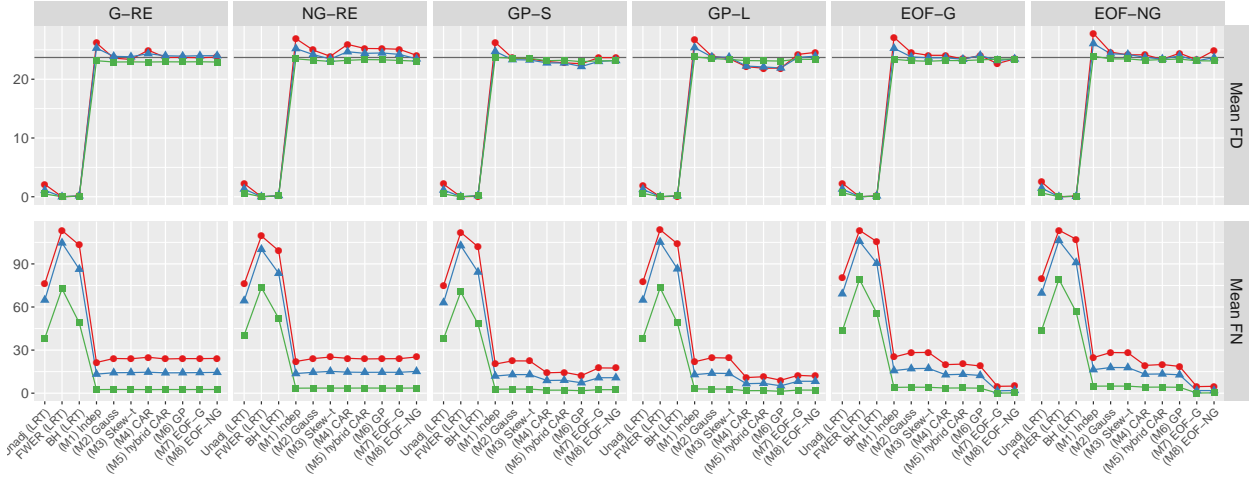


Figure A.5: Mean loss using the R_2 criteria, aggregated over the $N_{\text{rep}} = 100$ replicates, for schemes 1, 2, and 3. Note that the x -axis in each subgrid corresponds to the different methods/fitted models.

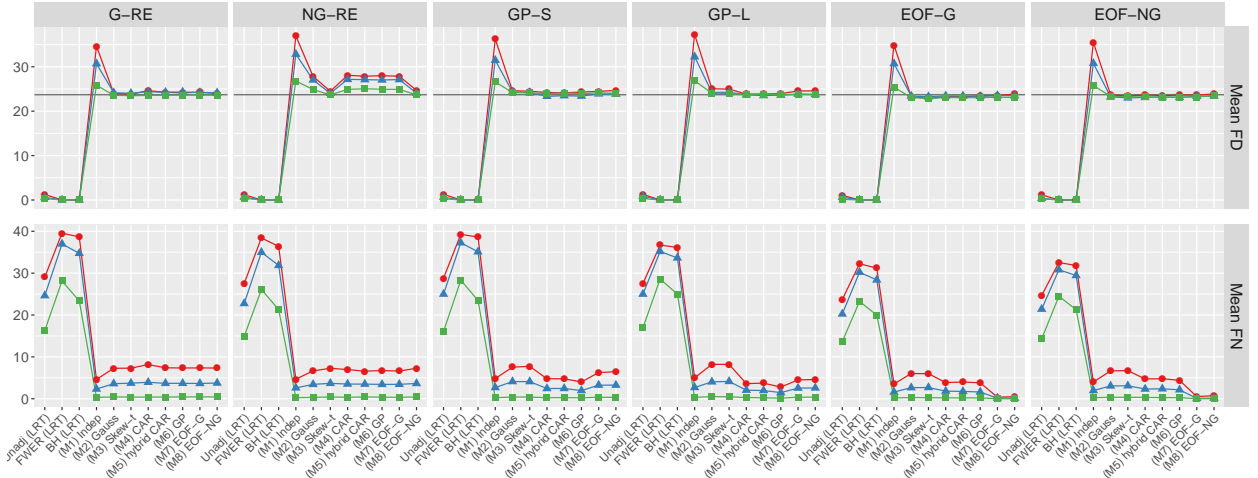
Scheme 1: approximately 85% true rejections



Scheme 2: approximately 50% true rejections



Scheme 3: approximately 15% true rejections



Ensemble size • 50 • 100 • 400

Figure A.6: Mean FD and FN using the R_3 criteria, aggregated over the $N_{\text{rep}} = 100$ replicates, for schemes 1, 2, and 3. Note that the x -axis in each subgrid corresponds to the different methods/fitted models. The target of $\gamma = 0.1M = 23.7$ is plotted for FD.

A.2 Results from simulation study with $M = 68$ regions

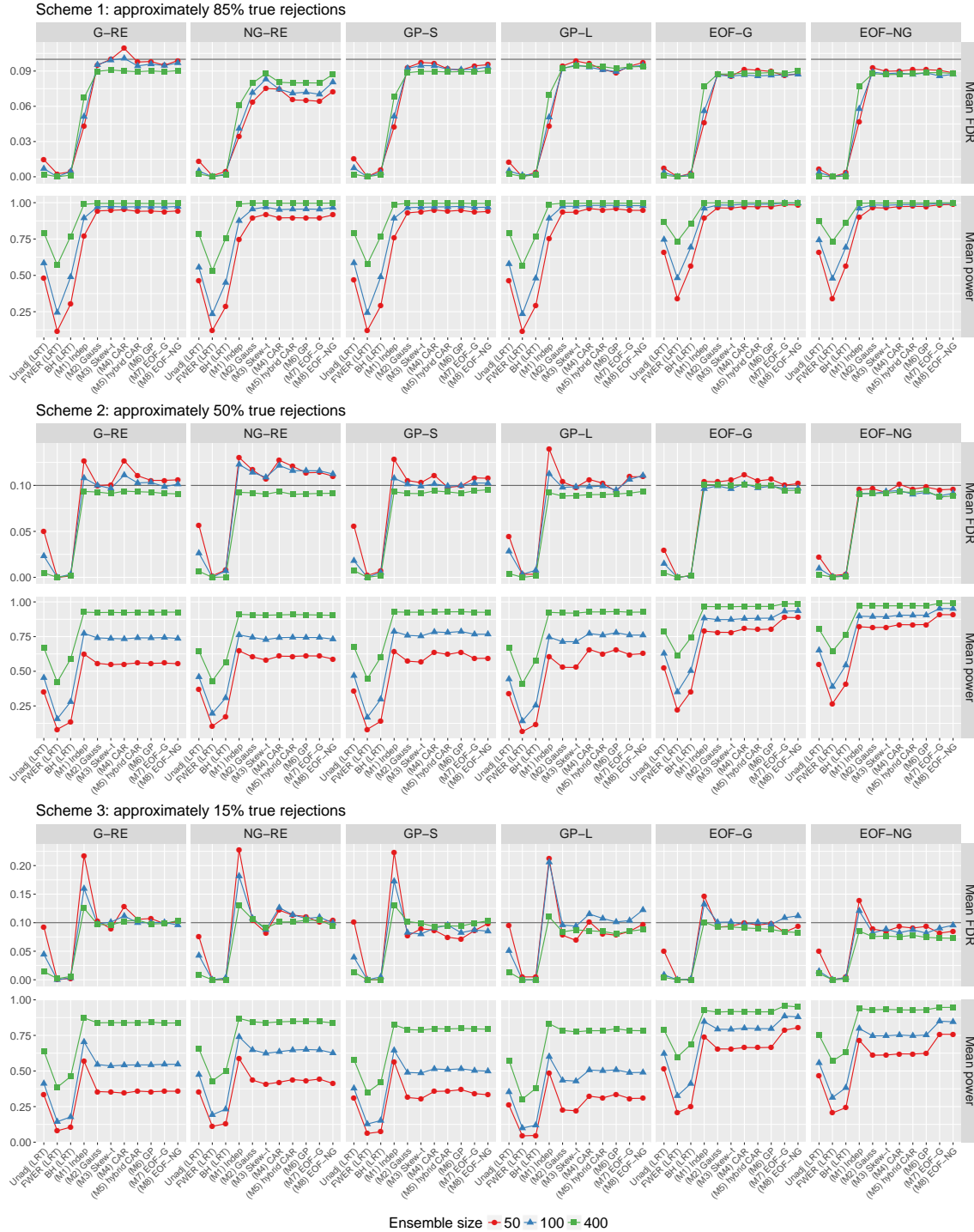


Figure A.7: Mean FDR and power using the R_1 criteria for the WRAF2 regions ($M = 68$), aggregated over the $N_{\text{rep}} = 100$ replicates, for schemes 1, 2, and 3. The target of $\alpha = 0.1$ is plotted for FDR.

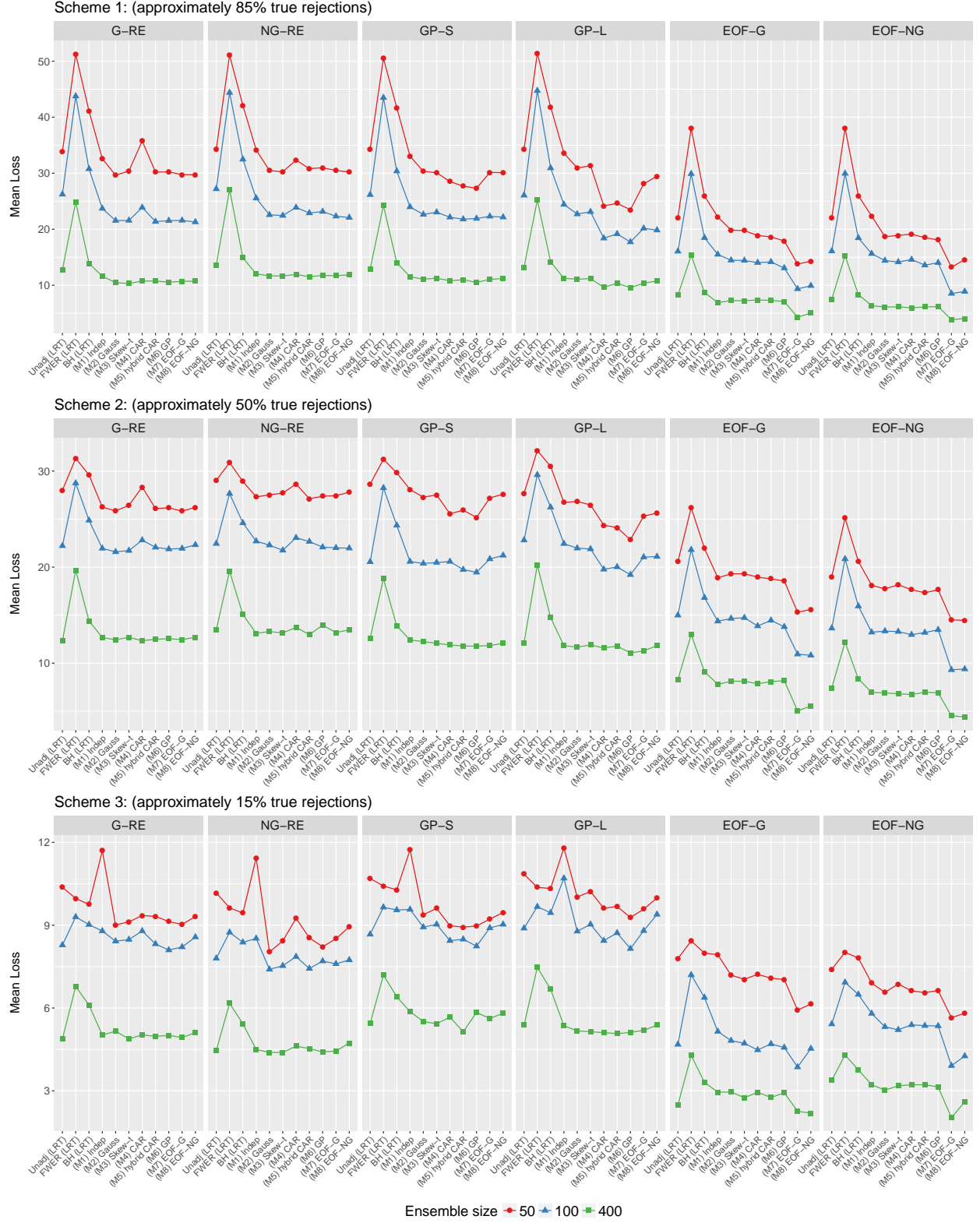


Figure A.8: Mean loss using the R_2 criteria for the WRAF2 regions ($M = 68$), aggregated over the $N_{\text{rep}} = 100$ replicates, for schemes 1, 2, and 3.

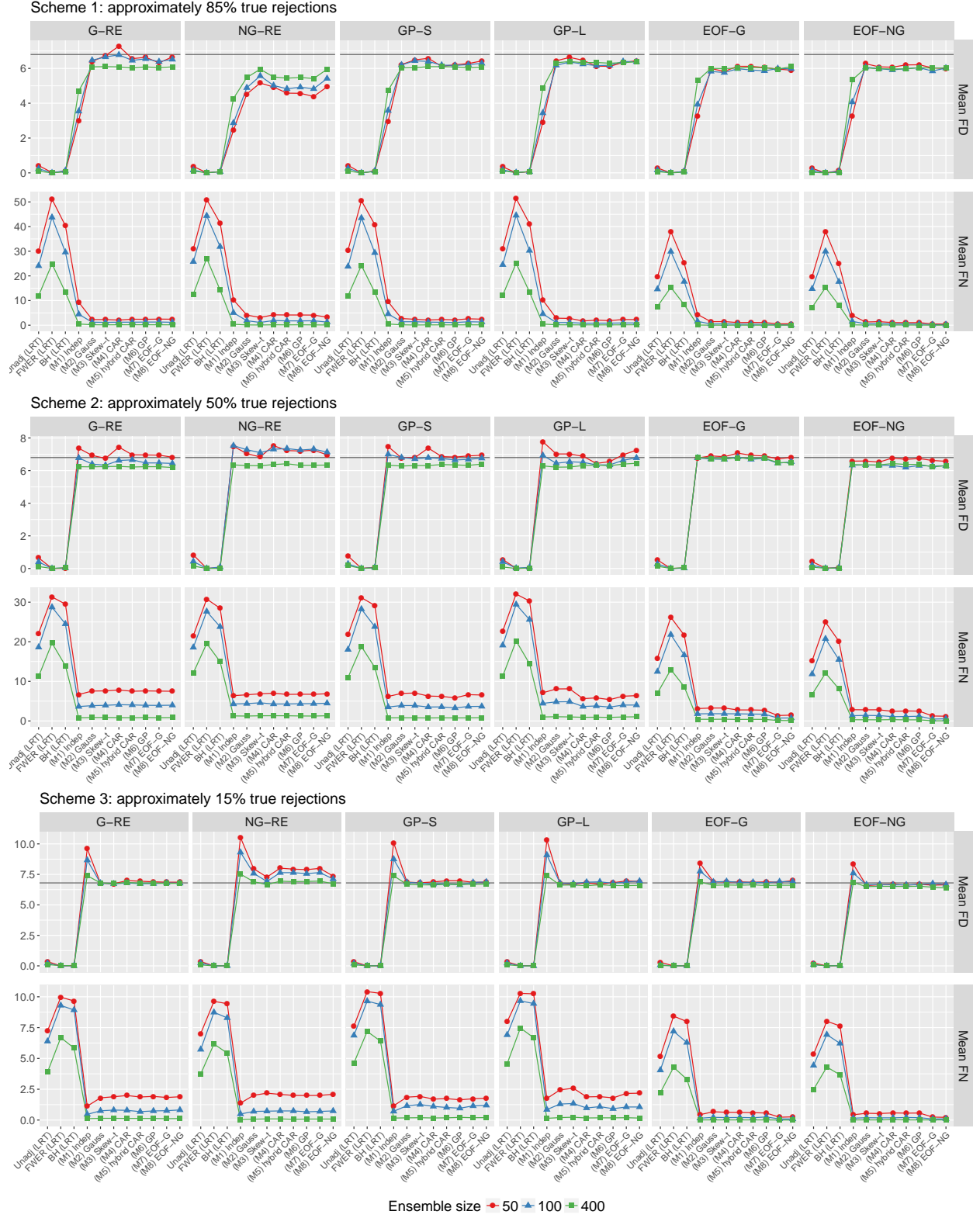


Figure A.9: Mean FD and FN using the R_3 criteria for the WRAF2 regions ($M = 68$), aggregated over the $N_{\text{rep}} = 100$ replicates, for schemes 1, 2, and 3. The target of $\gamma = 0.1M = 6.8$ is plotted for FD.

B Centered parameterization for the skew- t distribution

Note: the parameter symbols used in this section do not correspond to the symbols used in the main draft of the text.

[Azzalini and Capitanio \(2003\)](#) introduced the skew- t family of distributions, with probability density function

$$f_{ST}(y; \xi, \omega, \alpha, \nu) = \frac{2}{\omega} t_\nu \left(\frac{y - \xi}{\omega} \right) T_{\nu+1} \left(\frac{\alpha(y - \xi)}{\omega} \sqrt{\frac{\nu + 1}{\nu + (y - \xi)/\omega}} \right), \quad (\text{B.1})$$

where t_ν and T_ν denote the probability density and cumulative distribution function, respectively, of a standard t distribution with ν degrees of freedom. In (B.1), $\xi \in \mathbb{R}$ is a location parameter, $\omega \in \mathbb{R}^+$ is a scale parameter, $\alpha \in \mathbb{R}$ controls the skewness, and $\nu \in \mathbb{R}^+$ controls the tail behavior. Unfortunately, as noted by [Arellano-Valle and Azzalini \(2008\)](#) (and others), using the “direct” parameterization $\theta_D = (\xi, \omega, \alpha, \nu)$ has both theoretical and practical problems: for example, the likelihood behaves strangely for a neighborhood of $\alpha = 0$, in that the profile likelihood for α has a stationary point at 0. Furthermore, at $\alpha = 0$, the expected Fisher information is singular, even though all of the parameters are identifiable. In practical terms, this means that the parameter estimates (especially ξ and ω) can trade off with one another to give qualitatively similar results for an individual data set.

To address this problem, [Arellano-Valle and Azzalini \(2008\)](#) discuss a “centered” parameterization (for the skew-normal distribution; a corresponding result holds for the skew- t), originally introduced by [Azzalini and Capitanio \(2003\)](#). Instead of θ_D , the centered parameterization involves $\theta_C = (\mu, \sigma, \delta, \nu)$, where

$$\mu = \xi + \omega \sqrt{2/\pi} \frac{\alpha}{\sqrt{1 + \alpha^2}}, \quad -\infty < \mu < \infty,$$

$$\sigma = \omega \sqrt{1 - \frac{2}{\pi} \frac{\alpha^2}{1 + \alpha^2}}, \quad 0 < \sigma < \infty,$$

and

$$\delta = \frac{\alpha}{\sqrt{1 + \alpha^2}}, \quad -1 < \delta < 1,$$

with inverse transformations

$$\xi = \mu - \frac{\sigma}{\sqrt{1 - \frac{2}{\pi} \delta^2}} \sqrt{2/\pi} \delta, \quad \omega = \frac{\sigma}{\sqrt{1 - \frac{2}{\pi} \delta^2}}, \quad \alpha = \frac{\delta}{\sqrt{1 - \delta^2}}. \quad (\text{B.2})$$

(Note: ν is the same in both parameterizations.) Using θ_C avoids the problems associated with θ_D ; in practice, the likelihood associated with θ_C is given by (B.1), after substituting in (B.2).

C Prior specification for the parametric Bayesian models

In general, the priors used for all parameters will be proper but diffuse, with fixed hyperparameters. The details for each model are as follows; all of the priors below are for both $k \in \{F, C\}$.

M1 Beta-binomial, independent across regions

The only parameters in M1 are the probabilities themselves, which have already been assigned beta priors. The hyperparameters are set to $a_p = b_p = 1$, i.e., the probabilities are given an uniform prior.

M2 Gaussian random effects

The parameters in M2 are the scenario-specific mean μ_k and variance τ_k^2 , with priors

$$\mu_k \sim N(0, 10^2), \quad \tau_k^2 \sim U(0, 1000),$$

where $N(a, b)$ is the Gaussian distribution with mean a and variance b and $U(c, d)$ is the uniform distribution on the interval (c, d) .

M3 Skew- t random effects

Following [Arellano-Valle and Azzalini \(2008\)](#), M3 involves the scenario-specific “centered” parameters (see Appendix B) location μ_k , scale σ_k , skewness δ_k , and degrees of freedom ν_k . The prior distributions used are

$$\mu_k \sim N(0, 10^2), \quad \sigma_k \sim U(0, 100), \quad \delta_k \sim U(-1, 1), \quad 1/\nu_k \sim U(0, 1).$$

M4 CAR effects

The parameters in M4 are the scenario-specific mean μ_k and variance τ_k^2 ; however, because the CAR prior is improper, we fix $\mu_k = 0$ (see Appendix D). As before, $\tau_k^2 \sim U(0, 1000)$.

M5 Hybrid CAR/exchangeable effects

The parameters in M5 are the scenario-specific mean μ_k , variance τ_k^2 , and mixture parameter λ_k , with priors

$$\mu_k \sim N(0, 10^2), \quad \tau_k^2 \sim U(0, 1000), \quad \lambda_k \sim U(0, 1).$$

M6 Spatial Gaussian process effects

The parameters in M6 are the scenario-specific mean μ_k , variance τ_k^2 , and spatial “range” parameter ϕ_k , with priors

$$\mu_k \sim N(0, 10^2), \quad \tau_k^2 \sim U(0, 100), \quad \phi_k \sim U(0, c_\phi),$$

where $c_\phi = (1/2) \max\{\|\mathbf{s}_i - \mathbf{s}_j\|\}$, since the range of the Gaussian process would not be expected to exceed one-half of the maximum distance between the region centroids. Note that the smoothness parameter for the Matérn correlation function will be considered fixed, at 0.5 (corresponding to an exponential correlation function).

M7 EOF-based structure, Gaussian discrepancy

The parameters in M7 are the scenario-specific mean μ_k , EOF coefficients α_k , and variance τ_k^2 . As

before,

$$\mu_k \sim N(0, 10^2), \quad \tau_k^2 \sim U(0, 100^2);$$

because the basis functions are orthogonal, we a diagonal shrinkage prior for the EOF coefficients:

$$\boldsymbol{\alpha}_k \sim N(0, \sigma_\alpha^2 \mathbf{I}_p)$$

(i.e., $\Sigma_k^\alpha = \sigma_\alpha^2 \mathbf{I}_p$), where $\sigma_\alpha^2 \sim U(0, 100^2)$.

M8 EOF-based structure, skew- t discrepancy

The parameters in M8 are the scenario-specific mean μ_k , EOF coefficients $\boldsymbol{\alpha}_k$, scale σ_k , skewness δ_k , and degrees of freedom ν_k . All priors are as defined for M3 and M7.

D Markov chain Monte Carlo

The posterior distribution for each of the hierarchical models M2-M8 is not available in closed form, so we resort to Markov chain Monte Carlo (MCMC) methods to obtain samples from the joint posterior distribution for each model. All models are fit using the `nimble` software for R (NIMBLE Development Team, 2016). While the MCMC is straightforward for M2, M3, M5, M6, M7, and M8 (using standard Gibbs sampling with Metropolis Hastings steps), model M4 requires an adjustment to the standard MCMC (see the next section). The code used to fit these models are available in the online reproducibility documents.

D.1 Computational details for the CAR parameterization

Recall that computation for the CAR model is hindered by the fact that the CAR prior (9) is improper. This results in two problems: first, the random effects are identifiable only up to an additive constant; second, the CAR prior is undefined for the full random effects vector. While more sophisticated solutions to the first problem are possible, for the purposes of this work we simply set $\mu_k = 0$ to fix the identifiability problem.

Rue and Held (2005) outline steps to address the second problem. The CAR prior is

$$p(\boldsymbol{\beta}_k | \mathbf{Q}_k, \tau_k^2) \propto |\tau_k^{-2} \mathbf{Q}|^{1/2} \exp \left\{ -\frac{1}{2} \boldsymbol{\beta}_k^\top \mathbf{Q}_k \boldsymbol{\beta}_k \right\};$$

however, the rank of \mathbf{Q} is $M - 1$ ($\mathbf{1}^\top \mathbf{Q} = 0$), so the determinant $|\tau_k^{-2} \mathbf{Q}| = 0$. However, while the CAR prior is improper for an M -dimensional space, it is proper for a $(M - 1)$ -dimensional subspace. Following Rue and Held (2005), the prior contribution to the posterior is actually

$$\tilde{p}(\boldsymbol{\beta}_k | \mathbf{Q}_k, \tau_k^2) = (2\pi\tau_k^2)^{-\frac{(M-1)}{2}} \left(\prod_{i=1}^{M-1} \lambda_{ki} \right)^{1/2} \exp \left\{ -\frac{1}{2} \boldsymbol{\beta}_k^\top \mathbf{Q}_k \boldsymbol{\beta}_k \right\},$$

where $\{\lambda_{ki} : i = 1, \dots, M - 1\}$ are the non-zero eigenvalues of \mathbf{Q}_k .

E Further details on the simulation study

E.1 Simulation scheme for each true state

The six true states used as population distributions for the simulation study are listed in Table 2. The actual sampling procedure for each true state is now outlined.

First, for the Gaussian random effects (G-RE), the logit probabilities are simply draws from a Gaussian distribution:

$$\text{logit}(p_k) \stackrel{\text{iid}}{\sim} N(m_k, v_k^2).$$

Next, for the gamma random effects (NG-RE), the logit probabilities are draws from a shifted gamma distribution:

$$\text{logit}(p_k) \stackrel{\text{iid}}{\sim} G(a_k, b_k) - c_k,$$

where a_k and b_k are the shape and scale parameters, respectively. The Gaussian process samples (GP-S and GP-L) are first drawn collectively from

$$\text{logit}(\mathbf{p}_k) \sim N_M(m_k \mathbf{1}_M, \mathbf{S}),$$

where the elements of \mathbf{S} are $S_{ij} = v_k^2 \mathcal{M}_{g_k}(\|\mathbf{s}_i - \mathbf{s}_j\|/r_k)$ (where $\mathcal{M}_g(\cdot)$ is the Matérn correlation function and \mathbf{s}_i is the centroid of region i) and then centered to have an empirical mean of zero.

It is slightly less straightforward to generate samples from EOF-G and EOF-NG, especially because the generated data needs to have properties comparable to the other simulations (in terms of the correct proportion of true rejections and empirical variance of the true log risk ratio). The following (somewhat complicated) scheme made this possible (the k subscript has been omitted).

1. For $j = 1, \dots, p$ (where p is the total number of basis functions), calculate $a_j = B_j U_j$, where $B_j = -1I_j + 1(1 - I_j)$ with $I_j \sim \text{Ber}(0.5)$ and $U_j \sim U(l_j, u_j)$.
2. Draw $x_j \stackrel{\text{iid}}{\sim} N(0, v^2)$ (for EOF-G) or $x_j \stackrel{\text{iid}}{\sim} k[G(b, c) - d]$ (for EOF-NG).
3. Calculate the probabilities as $\mathbf{p} = \text{logit}^{-1} [m \mathbf{1}_M + \mathbf{H} \mathbf{a} + \mathbf{x}]$.

E.2 Fixed hyperparameter values for the true states

Tables E.1-E.5 contain the fixed hyperparameters used to sample draws from the fixed population distributions across the N_{rep} replicates. The values were determined after much trial and error, and were set according to two criteria: first, that the true proportion of rejections would match up with the corresponding scheme, and second, that the variance of the true log risk ratio (empirically, over many replicates) would be approximately 0.9.

Table E.1: Fixed hyperparameter values used for simulations from the Gaussian random effects (G-RE), across Schemes 1–3.

	Scheme 1	Scheme 2	Scheme 3
m_C	logit(0.08)	logit(0.08)	logit(0.08)
m_F	logit(0.03)	logit(0.08)	logit(0.19)
v_C^2, v_F^2	0.72^2	0.74^2	0.775^2

Table E.2: Fixed hyperparameter values used for simulations from the shifted gamma random effects (NG-RE), across Schemes 1–3. Note: a is the shape parameter and b is the scale parameter.

	Scheme 1	Scheme 2	Scheme 3
m_C	logit(0.08)	logit(0.08)	logit(0.08)
m_F	logit(0.03)	logit(0.08)	logit(0.18)
a_C, a_F	4	3.75	3.5
b_C, b_F	0.375	0.4	0.4286
c_C, c_F	1.5	1.5	1.5

Table E.3: Fixed hyperparameter values used for simulations from the spatial Gaussian process effects (GP-S and GP-L), across Schemes 1–3. Note: the distances in \mathbb{R}^3 are re-scaled to have a maximum of 1 unit.

	Scheme 1	Scheme 2	Scheme 3
m_C	logit(0.08)	logit(0.08)	logit(0.08)
m_F	logit(0.03)	logit(0.08)	logit(0.18)
v_C^2, v_F^2	0.6	0.6	0.6
r_C, r_F (short)	0.06	0.06	0.06
r_C, r_F (long)	0.10	0.10	0.10
g_C, g_F	2	2	2

Table E.4: Fixed hyperparameter values used for simulations from the EOF effects with Gaussian discrepancy (EOF-G), across Schemes 1–3.

	Scheme 1	Scheme 2	Scheme 3
m_C	logit(0.08)	logit(0.08)	logit(0.08)
m_F	logit(0.03)	logit(0.08)	logit(0.19)
$l_j, j = 1, \dots 5$	4.5	4.5	4.5
$u_j, j = 1, \dots 5$	5.5	5.5	5.5
$l_j, j = 6, \dots 10$	0.5	0.5	0.5
$u_j, j = 6, \dots 10$	2	2	2
$l_j, j = 11, \dots 20$	0	0	0
$u_j, j = 11, \dots 20$	0.1	0.1	0.1
v_C^2, v_F^2	0.01^2	0.01^2	0.01^2

Table E.5: Fixed hyperparameter values used for simulations from the EOF effects with gamma discrepancy (EOF-NG), across Schemes 1–3.

	Scheme 1	Scheme 2	Scheme 3
m_C	logit(0.08)	logit(0.08)	logit(0.08)
m_F	logit(0.03)	logit(0.08)	logit(0.19)
$l_j, j = 1, \dots 5$	4.5	4.5	4.5
$u_j, j = 1, \dots 5$	5.5	5.5	5.5
$l_j, j = 6, \dots 10$	0.5	0.5	0.5
$u_j, j = 6, \dots 10$	2	2	2
$l_j, j = 11, \dots 20$	0	0	0
$u_j, j = 11, \dots 20$	0.1	0.1	0.1
k_C, k_F	0.01	0.01	0.01
b_C, b_F	5	5	5
c_C, c_F	0.4	0.4	0.4
d_C, d_F	2	2	2